

#1

THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re the Application of : Naoki MATSUOKA, et al.

Filed : Concurrently herewith

For : PACKET SWITCH

Serial No. : Concurrently herewith



August 30, 2001

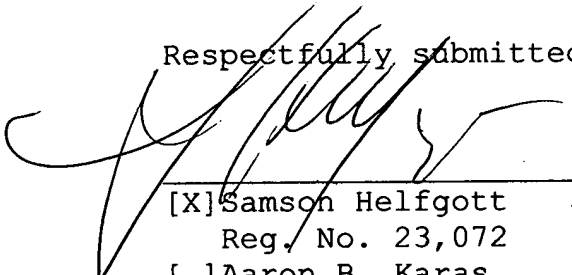
Assistant Commissioner of Patents
Washington, D.C. 20231

SUBMISSION OF PRIORITY DOCUMENT

S I R:

Attached herewith is Japanese Patent Application No. 2000-395741 of December 26, 2000 whose priority has been claimed in the present application.

Respectfully submitted



[X] Samson Helfgott
Reg. No. 23,072
[] Aaron B. Karas
Reg. No. 18,923

HELFGOTT & KARAS, P.C.
60th FLOOR
EMPIRE STATE BUILDING
NEW YORK, NY 10118
DOCKET NO.: FUJG 18.949
BHU:priority

Filed Via Express Mail
Rec. No.: EL639693967US
On: August 30, 2001
By: Brendy Lynn Belony
Any fee due as a result of this paper, not covered
by an enclosed check may be charged on Deposit Acct.
No. 08-1634.

日 本 国 特 許 庁
JAPAN PATENT OFFICE

J1000 U.S. PTO
09/942979
08/30/01

別紙添付の書類に記載されている事項は下記の出願書類に記載されて
いる事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed
with this Office

出 願 年 月 日

Date of Application:

2000年12月26日

出 願 番 号

Application Number:

特願2000-395741

出 願 人

Applicant(s):

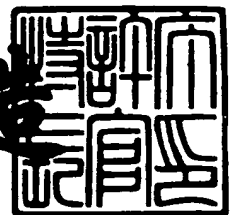
富士通株式会社

CERTIFIED COPY OF
PRIORITY DOCUMENT

2001年 6月 5日

特許庁長官
Commissioner,
Japan Patent Office

及川耕造



【書類名】 特許願

【整理番号】 0051384

【提出日】 平成12年12月26日

【あて先】 特許庁長官殿

【国際特許分類】 H04Q 3/00

【発明の名称】 パケットスイッチ

【請求項の数】 5

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 松岡 直樹

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 朝永 博

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 瓦井 健一

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 永田 将克

【特許出願人】

【識別番号】 000005223

【氏名又は名称】 富士通株式会社

【代理人】

【識別番号】 100103171

【弁理士】

【氏名又は名称】 雨貝 正彦

【電話番号】 03-3362-6791

【手数料の表示】

【予納台帳番号】 055491

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 0001848

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 パケットスイッチ

【特許請求の範囲】

【請求項 1】 N 本の入力回線のそれぞれに対応して設けられ、対応する前記入力回線を介して入力されるパケットを格納する N 個の入力バッファ部と、

それぞれが独立したスケジューリング処理を行うことにより、前記 N 個の入力バッファ部のそれぞれに格納された前記パケットの送出先となる M 本の出力回線のいずれかを決定する α 個のスケジューラ部と、

前記 N 個の入力バッファ部のそれぞれから出力される前記パケットを、前記スケジューラ部によって決定された送出先となる前記出力回線に出力するスイッチ部と、

を備え、前記 N 個の入力バッファ部において前記 α 個のスケジューラ部によるスケジューリング処理の結果を巡回的に使用することを特徴とするパケットスイッチ。

【請求項 2】 請求項 1 において、

前記スケジューラ部によるスケジューリング処理は、前記 N 個の入力バッファ部から送られてくるスケジューリング要求通知に対応して行われており、

前記 N 個の入力バッファ部のそれぞれは、前記スケジューリング要求通知の送出先となる前記スケジューラ部を分散させることを特徴とするパケットスイッチ。

【請求項 3】 請求項 2 において、

前記入力バッファ部は、前記 α 個のスケジューラ部のそれぞれに送出した前記スケジューリング要求通知の数を管理しており、この数が所定数に達した前記スケジューラ部に対する前記スケジューリング要求通知の送出動作を、この数が前記所定数より少なくなるまで遅らせることを特徴とするパケットスイッチ。

【請求項 4】 請求項 1 ～ 3 のいずれかにおいて、

前記スケジューラ部が前記スケジューリング処理に要する時間が、前記パケットの最小送出間隔の L 倍であるときに、

$L - \alpha$ が 1 以上に設定されており、

前記N個の入力バッファ部において前記 α 個のスケジューラ部の全てのスケジューリング処理の結果を巡回的に使用することを特徴とするパケットスイッチ。

【請求項5】 請求項1～4のいずれかにおいて、

前記スケジューラ部は、未使用回線を含む前記スケジューリング処理を行っており、実際に使用している前記入力回線および前記出力回線と前記未使用回線との間の読み替え処理を行うことにより、1回の前記スケジューリング処理によって複数のスケジューリング処理結果を得ることを特徴とするパケットスイッチ。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、入力回線を介して入力されたパケットを複数の出力回線に振り分けて出力するパケットスイッチに関する。

【0002】

【従来の技術】

一般家庭におけるインターネットユーザの増加やインターネットビジネスの成長を背景として、インターネットバックボーンネットワークの大容量化および高品質化が求められている。現在のインターネットは、品質保証のないベストエフォート型通信が主流であるが、次世代のインターネットでは、ベストエフォート型通信に加えて、音声やビデオ等のリアルタイムデータ通信の提供が期待されている。このため、ネットワークノードは、様々なトラフィックを柔軟に収容可能なテラビットクラスのスイッチング容量を有し、様々な通信サービスに適したサービス品質QoS (Quality of Service) を提供する必要がある。現在、テラビットクラスのスイッチング容量を有するノードを実現する手段としては、メモリアクセス速度の高速化が可能な入力バッファ型パケットスイッチが有望であると考えられている。

【0003】

ところで、入力バッファ型パケットスイッチにおけるスケジューリングアルゴリズムは、従来から様々なものが提案されている。これらをスケジューリング機能の配備方法に着目して整理すると、インタフェースカード等にスケジューリン

グ機能を分散して配備する方法と、専用のカード等にスケジューリング機能を集中して配備する方法とに大別することができる。

【 0 0 0 4 】

図 2 2 および図 2 3 は、従来の入力バッファ型パケットスイッチにおけるスケジューリング機能の配備状態を示す図である。図 2 2 にはスケジューリング機能を分散して配備したパケットスイッチの構成が示されている。図 2 3 にはスケジューリング機能を集中して配備したパケットスイッチの構成が示されている。

【 0 0 0 5 】

図 2 2 に示すように、スケジューリング機能を分散配備した入力バッファ型パケットスイッチには、入力バッファ（B U F）とともにスケジューラが備わったインタフェースカード 1 0 0 が入力回線の数だけ設けられている。隣接するインタフェースカード 1 0 0 内のスケジューラ同士を相互接続することにより、各スケジューラ間でスケジューリング情報の送受信（未確定回線の通知）が可能になり、全ての入力回線のスケジューリングが実施される。入力回線数が増減した場合には、入力回線数に合わせてインタフェースカード 1 0 0 の数も増減させればよい。図 2 2 に示した構成を有する入力バッファ型パケットスイッチは、拡張性に優れるという利点を有する。

【 0 0 0 6 】

一方、図 2 3 に示すように、スケジューリング機能を集中配備した入力バッファ型パケットスイッチは、全ての入力回線のスケジューラ機能部（S C H）を一つのスケジューラカード 1 1 0 に集約しており、これによって全ての入力回線のスケジューリングが実施される。各スケジューラ機能部間を接続する配線の距離が短いため、この配線により生じる遅延が少なく、信号の遅延時間に起因する実装時の制約が少ないという利点がある。

【 0 0 0 7 】

【発明が解決しようとする課題】

ところで、スケジューリング機能を分散配備した従来の入力バッファ型パケットスイッチは、実装された各インタフェースカード 1 0 0 に備わったスケジューラ間を相互接続する必要があるため、接続線の配線長による遅延量が多くなり、

高速なパケットスイッチでは、実装時の制約が多いという問題があった。

【0008】

また、スケジューリング機能部を集中配備した従来の入力バッファ型パケットスイッチは、実際に収容される入力回線数が少ない場合であっても、スケジューラカード110には、常に最大数のスケジューラ機能部を備えておく必要があり、しかも実装後はその数が固定であるため、無駄が多いとともに拡張性に欠けるという問題があった。

【0009】

本発明は、このような点に鑑みて創作されたものであり、その目的は、実装時の制約が少なく、無駄な構成を低減でき、しかも拡張性を有するパケットスイッチを提供することにある。

【0010】

【課題を解決するための手段】

上述した課題を解決するために、本発明のパケットスイッチは、N個の入力バッファ部と α 個のスケジューラ部とスイッチ部とを備えている。N個の入力バッファ部は、N本の入力回線のそれぞれに対応して設けられ、対応する入力回線を介して入力されるパケットを格納する。 α 個のスケジューラ部は、それぞれが独立したスケジューリング処理を行うことにより、N個の入力バッファ部のそれぞれに格納されたパケットの送出先となるM本の出力回線のいずれかを決定する。スイッチ部は、N個の入力バッファ部のそれぞれから出力されるパケットを、スケジューラ部によって決定された送出先となる出力回線に出力する。上述したN個の入力バッファ部において α 個のスケジューラ部によるスケジューリング処理の結果を巡回的に使用する。複数のスケジューラ部は互いに独立してスケジューリング処理を行っているため、これらの間の競合制御等を行う必要がなく、このための信号線による信号の遅延等の問題もないことから、実装時の制約が大幅に緩和される。また、スケジューラ部は必要な数だけ用意すればよいため無駄な構成がなく、しかも後に追加することも可能であるため、拡張性に富んだパケットスイッチを実現することができる。

【0011】

特に、上述したスケジューラ部によるスケジューリング処理は、N個の入力バッファ部から送られてくるスケジューリング要求通知に対応して行われており、N個の入力バッファ部のそれぞれは、スケジューリング要求通知の送出先となるスケジューラ部を分散させることが望ましい。スケジューリング要求通知の送出先となるスケジューラ部を分散させることにより、処理の負担がスケジューラ部間で偏ることを防止することができる。

【 0 0 1 2 】

また、上述した入力バッファ部は、M本の出力回線のそれぞれを送出先とするパケットを格納するM個のキューを有しており、これらM個のキューのそれぞれ毎にスケジューリング要求通知の送出先となるスケジューラ部を巡回させることが望ましい。スケジューリング要求通知の送出先となるスケジューラ部をキュー毎に巡回させることにより、各スケジューラ部に対するスケジューリング要求通知を分散させることができる。

【 0 0 1 3 】

また、上述した入力バッファ部は、入力回線毎にスケジューリング要求通知の送出先となるスケジューラ部を巡回させることが望ましい。スケジューリング要求通知の送出先となるスケジューラ部を入力回線毎に巡回させることにより、各スケジューラ部に対するスケジューリング要求通知を分散させることができる。

【 0 0 1 4 】

上述した入力バッファ部は、単位時間毎にスケジューリング要求通知の送出先となるスケジューラ部を巡回させることが望ましい。スケジューリング要求通知の送出先となるスケジューラ部を単位時間毎に巡回させることにより、各スケジューラ部に対するスケジューリング要求通知を分散させることができる。

【 0 0 1 5 】

また、上述した入力バッファ部は、 α 個のスケジューラ部のそれぞれについて未処理のスケジューリング要求通知の数を調べ、この数が少ないスケジューラ部に次のスケジューリング要求通知を送ることが望ましい。実際に送出されたスケジューリング要求通知の数が少ないスケジューラ部に優先的にスケジューリング要求通知を送ることにより、各スケジューラ部に対するスケジューリング要求通

知を均等に分散させることができる。

【 0 0 1 6 】

また、上述した入力バッファ部は、 α 個のスケジューラ部のそれぞれに送出したスケジューリング要求通知の数を管理しており、この数が所定数に達したスケジューラ部に対するスケジューリング要求通知の送出動作を、この数が所定数より少なくなるまで遅らせることが望ましい。入力バッファ部にスケジューリング要求通知の数を管理する機能を分散させることにより、各スケジューラ部の処理負担を軽減することができる。

【 0 0 1 7 】

また、上述したスケジューラ部がスケジューリング処理に要する時間が、パケットの最小送出間隔の L 倍であるときに、スケジューラ部の数 α は、倍数 L 以上の値に設定されていることが望ましい。スケジューラ部の数 α を倍数 L 以上に設定することにより、 α 個のスケジューラ部全体として遅延なくスケジューリング処理を行うことができる。

【 0 0 1 8 】

また、 $L - \alpha$ が 1 以上に設定されており（すなわち、 L 個より多いスケジューラが備わっており）、 N 個の入力バッファ部において α 個のスケジューラ部の全てのスケジューリング処理の結果を巡回的に使用することが望ましい。処理の遅延を来さないために必要な数よりも多い数のスケジューラ部を備えることにより、その一部に障害が発生した場合であっても、残りのスケジューラ部によって正常にスケジューリング処理を継続することができる。

【 0 0 1 9 】

また、 $L - \alpha$ が 1 以上に設定されている場合に、 $\alpha - L$ 個のスケジューラ部を冗長系として用いるとともに、この冗長系以外のスケジューラ部に障害が発生したときに、代わりに冗長系のスケジューラ部を用いることが望ましい。障害発生時に冗長系のスケジューラ部を用いてスケジューリング処理を行うことにより、遅延を来すことなく正常なスケジューリング処理を継続することができる。

【 0 0 2 0 】

また、入力回線の数 N と出力回線の数 M に応じて、スケジューラ部の数 α およ

びスケジューリング処理の時間を可変に設定することが望ましい。回線数が少ない場合にはスケジューラ部の数や処理時間を少なく設定し、反対に回線数が多い場合にはスケジューラ部の数や処理時間を多く設定すればよい。ため、スイッチ部に収容する回線数等に応じて無駄のない最適な構成を実現することができる。

【 0 0 2 1 】

また、上述したスケジューラ部は、未使用回線を含むスケジューリング処理を行っており、実際に使用している入力回線および出力回線と未使用回線との間の読み替え処理を行うことにより、1回のスケジューリング処理によって複数のスケジューリング処理結果を得ることが望ましい。スケジューラ部の数や処理時間を変更することなく、回線数の変更に対応することができ、拡張性に富んだパケットスイッチを実現することができる。

【 0 0 2 2 】

【発明の実施の形態】

以下、本発明を適用した一実施形態のパケットスイッチについて詳細に説明する。

図1は、本発明を適用した一実施形態のパケットスイッチの構成を示す図である。図1に示すように、本実施形態のパケットスイッチ100は、スイッチ部(SW)10、N個の入力バッファ部20、 α 個のスケジューラ部30を含んで構成されている。

【 0 0 2 3 】

スイッチ部10は、N本の入力回線のいずれかから入力されたパケットをM本の出力回線のいずれかに出力する。本実施形態では、このスイッチ部10は、内部にパケットを保持することができないバッファレスタイプのものが用いられている。N個の入力バッファ部20のそれぞれは、スイッチ部10に収容されたN本の入力回線のそれぞれに対応して設けられている。各入力バッファ部20は、スイッチ部10に収容されたM本の出力回線に対応したM個の論理キューVOQ (Virtual Output Queue) を有しており、1対1に対応する入力回線を介して入力された到着パケットを、出力先となる出力回線に対応する論理キューVOQに蓄積するとともに、いずれかのスケジューラ部30から出力許可を示すグラント

通知が送られてきたときに、このグラント通知で指定された論理キューVOQの先頭パケットを読み出す。 α 個のスケジューラ部30のそれぞれは、各入力回線を介して入力されたパケットの送出先となる出力回線を決定するスケジューリング処理を行う。このようにして α 個のスケジューラ部30のそれぞれにおいてスケジューリング処理を並行して行うことにより、各スケジューラ部30による処理負担を軽減することができる。本明細書では、このようにして複数のスケジューラ部30によって並行して行うスケジューリング処理を「負荷分散スケジューリング処理」と称する。

【0024】

図2は、入力バッファ部20の詳細構成を示す図である。図2に示すように、入力バッファ部20は、上述したM個の論理キューVOQを含むパケットバッファ22の他に、要求振り分け部24および読み出し指示部26を有している。要求振り分け部24は、パケットバッファ22内の各論理キューVOQに蓄積されたパケットに対応するスケジューリング要求通知であるリクエスト通知を、 α 個のスケジューラ部30のそれぞれに均等に分散させて送出する。読み出し指示部26は、 α 個のスケジューラ部30から送られてくるグラント通知の内容にしたがって、パケットの送出タイミングが到来したパケットバッファ22内の論理キューVOQに読み出し指示を与える。

【0025】

図3は、スケジューラ部30の詳細構成を示す図である。図3に示すように、スケジューラ部30は、要求数管理部32およびスケジューリング制御部34を有している。要求数管理部32は、入力回線と論理キューVOQの組合せに対応するリクエスト通知が送られてきたときに、この組合せに対応するリクエスト数をインクリメントし、このリクエスト数をスケジューリング確定時にデクリメントする。スケジューリング制御部34は、所定のスケジューリングアルゴリズムにしたがって、ある時刻Tに送出する各入力回線毎のパケットの出力先回線を決定する。具体的には、スケジューリング制御部34は、入力回線のそれぞれについて、要求数管理部32によってカウントされているリクエスト数が1以上の出力回線の中からいずれかを選択する。

【 0 0 2 6 】

本実施形態では、 α 個のスケジューラ部 3 0 によって並列にスケジューリング処理が行われるため、1 個のスケジューラ部によってスケジューリング処理を行う場合に比べて、各スケジューラ部 3 0 は、 α 倍の時間をかけてスケジューリング処理を行うことができる。例えば、各入力回線を介して入力されたパケットを、対応する出力回線に送出するために必要な時間を「1 パケット時間」とすると、各スケジューラ部 3 0 は、N 本の入力回線のそれぞれを介して入力されたパケットの送出先となる出力回線を決定するスケジューリング処理を α パケット時間以内に行えばよい。

【 0 0 2 7 】

本実施形態のパケットスイッチ 1 0 0 はこのような構成を有しており、次にその動作をスケジューリング処理に着目して説明する。

(1) 入力バッファ部 2 0 の動作

(1 a) 対応する入力回線を介してパケットが入力された場合に、パケットバッファ 2 2 は、このパケットの送出先となる出力回線に対応する論理キュー VOQ にこのパケットを格納する。

【 0 0 2 8 】

(1 b) パケットバッファ 2 2 内のいずれかの論理キュー VOQ にパケットが格納されると、要求振り分け部 2 4 は、このパケットに対応するスケジューリング処理を要求するリクエスト通知をいずれかのスケジューラ部 3 0 に送る。このとき、リクエスト通知の送信先となるスケジューラ部 3 0 を分散させる。

【 0 0 2 9 】

(1 c) いずれかのスケジューラ部 3 0 から送信許可を示すグラント通知が送られてくると、読み出し指示部 2 6 は、このグラント通知によって指定された論理キュー VOQ の先頭のパケットを、この論理キュー VOQ に 1 対 1 に対応する出力回線に向けて出力する動作をパケットバッファ 2 2 に指示する。

【 0 0 3 0 】

(2) スケジューラ部 3 0 の動作

α 個のスケジューラ部 3 0 の動作は互いに独立であり、各入力バッファ部 2 0

から送られてくるリクエスト通知の数（リクエスト数）を論理キューVOQ毎に管理する。したがって、各スケジューラ部30は、自分宛てに送られてきたリクエスト通知の有無を基に、スケジューリング処理を実施する。また、スケジューリングアルゴリズムは、従来から用いられている各種の手法を採用することができる。以下では、ラウンドロビンを用いて、各出力回線毎に一つの論理キューVOQを選択する場合の動作を説明する。

【0031】

(2a) 各入力バッファ部20からリクエスト通知が送られてくると、要求数管理部32は、リクエスト通知の内容を調べ、入力回線と論理キューVOQの組合せ毎にリクエスト数をカウントする。

(2b) スケジューリング制御部34は、一の入力回線に着目し、送られてきたリクエスト通知によって指定されており、しかもそれまでのスケジューリング処理によって選択されていない（未確定の）論理キューVOQの中から一つを選択する。この入力回線毎に論理キューVOQを選択する処理は、全ての入力回線について一巡するまで順番に行われる。これにより、全ての入力回線のそれぞれに入力されたパケットの出力先となる出力回線が決定される。

【0032】

(2c) スケジューリング制御部34は、入力回線毎に選択した論理キューVOQを指定したグラント通知を、各入力バッファ部20内の読み出し指示部26に送る。

具体例1

図4は、負荷分散スケジューリング処理の具体例を示す図である。同図には、説明を簡略化するために、入力回線数、出力回線数、スケジューラ部30の数 α が全て2の場合に対応したスケジューリング処理の具体例が示されている。以下の説明および図4においては、2本の入力回線、2本の出力回線、各出力回線に対応した論理キューVOQのそれぞれに通し番号#0、#1が付されており、これらの識別が行われている。また、2つのスケジューラ部30の一方をスケジューラ部#0（SCH#0）、他方をスケジューラ部#1（SCH#1）とする。同様に、入力回線#0、#1のそれぞれに対応する2つの入力バッファ部20の

一方を入力バッファ部 # 0、他方を入力バッファ部 # 1 とする。また、これらに含まれる各構成についても、「# 0」あるいは「# 1」を名称の後に付して、対応する入力回線を区別するものとする。

【0033】

入力バッファ部の動作

一方の入力バッファ部 # 0 には、入力回線 # 0 を介して 4 つのパケット (1)、(2)、(3)、(4) が順番に入力される。この中で、第 1、第 2、第 3 パケット (1)、(2)、(3) は送出先として出力回線 # 0 が指定されており、第 4 パケット (4) のみが送出先として出力回線 # 1 が指定されているものとする。したがって、入力バッファ部 # 0 内のパケットバッファ # 0 は、第 1、第 2、第 3 パケット (1)、(2)、(3) を論理キュー VOQ # 0 に格納し、第 4 パケット (4) を論理キュー VOQ # 1 に格納する。要求振り分け部 # 0 は、第 1、第 3、第 4 のパケット (1)、(3)、(4) が到着したときに、一方のスケジューラ部 # 0 にリクエスト通知を送るとともに、第 2 のパケット (2) が到着したときに、他方のスケジューラ部 # 1 にリクエスト通知を送る。このリクエスト通知には、入力回線番号、出力回線番号 (論理キュー VOQ の番号)、イネーブル等が含まれている。

【0034】

また、他方の入力バッファ部 # 1 には、入力回線 # 1 を介して第 5、第 6 のパケット (5)、(6) が順番に入力される。これら 2 つのパケット (5)、(6) は、両方とも送出先として出力回線 # 1 が指定されているものとする。したがって、入力バッファ部 # 1 内のパケットバッファ # 1 は、これら第 5 および第 6 のパケット (5)、(6) を論理キュー VOQ # 1 に格納する。要求振り分け部 # 1 は、第 5 のパケット (5) が到着したときに、一方のスケジューラ部 # 0 にリクエスト通知を送るとともに、第 6 のパケット (6) が到着したときに、他方のスケジューラ部 # 1 にリクエスト通知を送る。

【0035】

スケジューラ部の動作

スケジューラ部 # 0、# 1 は、独立して動作しており、それぞれに含まれる要

求数管理部 # 0、# 1 によって、受信したリクエスト通知の数を論理キュー V O Q 毎に管理する。また、スケジューラ部 # 0、# 1 は、それぞれの要求管理部 # 0、# 1 によって管理されるリクエスト数に基づいて、以下に示す負荷分散スケジューリング処理を実施する。なお、ラウンドロビンを用いたスケジューリング処理が、入力回線 # 1、入力回線 # 0 の順番で行われるものとする。

【0036】

例えば、スケジューラ部 # 1 内のスケジューリング制御部 # 1 は、まず入力回線 # 1 に対応する全ての論理キュー V O Q についてリクエスト通知の有無を判定し、リクエスト通知が存在し、かつパケットの送出先が未確定の論理キュー V O Qの中から、ラウンドロビンを用いて一つを選択する。図 4 に示した例では、入力回線 # 1 (i # 1) については論理キュー V O Q # 1 のみに対応するリクエスト通知が存在し、しかもこの論理キュー V O Q # 1 は未確定状態にあるため、論理キュー V O Q # 1 が選択される。

【0037】

次に、スケジューラ部 # 1 内のスケジューリング制御部 # 1 は、上述した処理によって選択された論理キュー V O Q # 1 を確定済みとし、入力回線番号を変えて次の入力回線 # 0 について同様のスケジューリング処理を実施する。図 4 に示した例では、入力回線 # 0 については論理キュー V O Q # 0、# 1 のそれぞれに対応するリクエスト通知が存在し、この中で論理キュー V O Q # 1 は確定済みであるため、論理キュー V O Q # 0 が選択される。

【0038】

このようにして、スケジューラ部 # 1 内のスケジューリング制御部 # 1 は、 α ($= 2$) パケット時間内で、全ての入力回線分のスケジューリング処理を行う。

同様にして、スケジューラ部 # 0 内のスケジューリング制御部 # 0 においても、 α パケット時間内で、全ての入力回線分のスケジューリング処理が行われる。その結果、ある周期 A におけるスケジューラ # 0、# 1 のスケジューリング結果が以下のように決定される。

【0039】

スケジューラ # 0 : 入力回線 # 0 → 出力回線 # 1

スケジューラ # 1 : 入力回線 # 0 → 出力回線 # 0

入力回線 # 1 → 出力回線 # 1

これらの α 個 (2 個) の各スケジューラ部 # 0、# 1 の各スケジューリング結果は、全ての入力バッファ部 # 0、# 1 に対してグラント通知として送られる。入力バッファ部 # 0、# 1 内の読み出し指示部 # 0、# 1 は、1 パケット時間毎に異なるスケジューラ部 # 0、# 1 のスケジューリング結果を巡回的に使用し、パケットの送出先として指定された出力回線に対応する論理キュー VOQ の先頭からパケットを読み出す指示をパケットバッファ # 0、# 1 に送る。

【0040】

図4に示した例では、上述したスケジューリング処理が行われた周期Aの次の周期Bにおいて、最初の1パケット時間 ($T = a$) では、スケジューラ部 # 0 のスケジューリング結果が使用される。すなわち、入力バッファ部 # 0 内の論理キュー VOQ # 1 の先頭パケット (第4のパケット (4)) が出力回線 # 1 に向けて送出される。

【0041】

また、周期Bに含まれる次の1パケット時間 ($T = a + 1$) では、スケジューラ部 # 1 のスケジューリング結果が使用される。すなわち、入力バッファ部 # 0 内の論理キュー VOQ # 0 の先頭パケット (第1のパケット (1)) が出力回線 # 0 に向けて送出されるとともに、入力バッファ部 # 1 内の論理キュー VOQ # 1 の先頭パケット (第5のパケット (5)) が出力回線 # 1 に向けて送出される。

【0042】

ここで、各入力バッファ部 # 0、# 1 から送られるリクエスト通知と、実際に送出されるパケットとの関係に着目すると、1対1になっていないことがわかる。例えば、入力バッファ部 # 0 に第1のパケット (1) が到着して論理キュー VOQ # 0 に格納された際にリクエスト通知 (1) がスケジューラ部 # 0 に送られるが、この第1のパケット (1) が実際に出力回線 # 0 に向けて送出されるのは、周期Bに含まれる後半の1パケット時間 ($T = a + 1$) のタイミングであり、しかもスケジューラ部 # 1 によるスケジューリング結果に基づいて送出される。

これは、リクエスト通知を各スケジューラ部 # 0、# 1 に分散させたことによる同一論理キュー V O Q 内でのパケットの追い越しの発生や、H O L (Head of Line) ブロッキング等の影響を回避するために、リクエスト通知とグラント通知との対応を考えず、送られてきたグラント通知に対して常に論理キュー V O Q の先頭のパケットを読み出すようにしているためである。

【 0 0 4 3 】

具体例 2

図 5 は、負荷分散スケジューリング処理の他の具体例を示す図である。例えば、入力回線数、出力回線数、スケジューラ部 3 0 の数 α が全て 4 の場合に対応したスケジューリング処理の具体例が示されている。また、図 5 において、S C H # 0 ~ S C H # 3 は 4 つのスケジューラ部 3 0 (スケジューラ部 # 0 ~ # 3) のそれぞれを示している。また、「Scheduling → Result#0-1」は、スケジューリング処理によって「# 0 - 1」で特定されるスケジューリング結果が得られたことを示している。

【 0 0 4 4 】

図 5 に示す具体例では、4 つのスケジューラ部 # 0 ~ # 3 が独立にスケジューリング処理を行っており、4 パケット時間に相当する第 1 の周期 ($T = 0 \sim 3$) の最後に 4 つのスケジューリング結果 (Result#0-1 ~ #3-1) が得られる。これらのスケジューリング結果に基づいて、次の 4 パケット時間に相当する第 2 の周期 ($T = 4 \sim 7$) におけるパケット送出が行われる。

【 0 0 4 5 】

具体的には、第 2 の周期の最初の 1 パケット時間 ($T = 4$) において、スケジューラ部 # 0 (S C H # 0) によるスケジューリング結果「Result#0-1」が用いられる。この結果、以下に示す内容でパケットの送出が実施される。なお、左の項は、パケットが入力される入力回線を示しており、右の項は、スケジューリング結果に基づいてこのパケットがどの出力回線に送出されたかを示している。

【 0 0 4 6 】

- ・ 入力回線 # 0 (i # 0) → 出力回線 # 2 (o # 2)
- ・ 入力回線 # 1 (i # 1) → 出力回線 # 0 (o # 0)

・入力回線 # 2 (i # 2) → 出力回線 # 3 (o # 3)

・入力回線 # 3 (i # 3) → 出力回線 # 1 (o # 1)

また、第 2 の周期の次の 1 パケット時間 (T = 5) において、スケジューラ部 # 1 (S C H # 1) によるスケジューリング結果「Result#1-1」が用いられる。
この結果、以下に示す内容でパケットの送出が実施される。

【 0 0 4 7 】

・入力回線 # 0 (i # 0) → 出力回線 # 2 (o # 2)

・入力回線 # 1 (i # 1) → (未送出)

・入力回線 # 2 (i # 2) → 出力回線 # 0 (o # 0)

・入力回線 # 3 (i # 3) → 出力回線 # 3 (o # 3)

第 2 の周期の次の 1 パケット時間 (T = 6) において、スケジューラ部 # 2 (S C H # 2) によるスケジューリング結果「Result#2-1」が用いられる。この結果、以下に示す内容でパケットの送出が実施される。

【 0 0 4 8 】

・入力回線 # 0 (i # 0) → (未送出)

・入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

・入力回線 # 2 (i # 2) → 出力回線 # 0 (o # 0)

・入力回線 # 3 (i # 3) → 出力回線 # 2 (o # 2)

第 2 の周期の最後の 1 パケット時間 (T = 7) において、スケジューラ部 # 3 (S C H # 3) によるスケジューリング結果「Result#3-1」が用いられる。この結果、以下に示す内容でパケットの送出が実施される。

【 0 0 4 9 】

・入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

・入力回線 # 1 (i # 1) → 出力回線 # 2 (o # 2)

・入力回線 # 2 (i # 2) → 出力回線 # 1 (o # 1)

・入力回線 # 3 (i # 3) → 出力回線 # 3 (o # 3)

このように、4 つのスケジューラ部 # 0 ~ # 3 のそれぞれが独立にスケジューリング処理を行うことにより、互いに競合する内容のスケジューリング結果が得られても、各スケジューラ部 # 0 ~ # 3 によるスケジューリング結果が用いられ

るタイミングが異なるため、特に不都合は生じない。このため、各スケジューラ部 #0～#3 間の競合制御を行う接続線等が不要であり、この接続線によって生じる信号の遅延がなくなるため、実装時の制約が大幅に緩和される。また、スケジューラ部は必要な数だけ用意すればよいため無駄な構成がなく、しかも後に追加することも可能であるため、拡張性に富んだパケットスイッチ 100 を実現することができる。

【0050】

なお、上述した説明では、 α 個のスケジューラ部 #0、…、# $\alpha-1$ のそれぞれにおいて得られたスケジューリング結果を各入力バッファ部に同時に送り、各入力バッファ部において、1 パケット時間毎に異なるスケジューラ部によるスケジューリング結果を巡回的に用いるようにしたが、スケジューラ部からスケジューリング結果を送るタイミングを各スケジューラ部毎に 1 パケット時間ずらして順番に各入力バッファ部に送るようにしてもよい。

【0051】

リクエスト通知の分散

ところで、上述したように α 個のスケジューラ部 30 を用いて負荷分散スケジューリング処理を行うためには、各入力バッファ部 20 から送出するリクエスト通知を全てのスケジューラ部 30 に均等に分散させることが望ましい。

【0052】

次に、入力バッファ部 20 からスケジューラ部 30 に送るリクエスト通知の分散方法について説明する。

論理キュー VOQ 毎にリクエスト通知を分散させる場合

図 6 は、入力バッファ部 20 内の論理キュー VOQ 毎にリクエスト通知を分散させる場合の具体例を示す図である。なお、図 6 および以下に示す図 7～図 9 においては、4 つの入力バッファ部 20 (#0～#3) および 4 つのスケジューラ部 30 (#0～#3) が備わっているものとする。

【0053】

各入力バッファ部 #0～#3 内の要求振り分け部 #0～#3 は、論理キュー VOQ 毎の振り分けポインタを有しており、入力回線を介してパケットが到着した

ときに、この振り分けポインタの値が示すスケジューラ部に対してリクエスト通知を送る。図6では入力バッファ部#0、#3内の丸付き数字によってポインタ値が示されている。

【0054】

以下では、入力回線#0に対応する入力バッファ部#0内の要求振り分け部#0に着目して具体的な動作を説明する。なお、初期状態において、入力回線#0に対応する4つの論理キュー#0～#3の各振り分けポインタの値が0、1、3、1に設定されているものとする。

【0055】

T=0において、送出先として出力回線#3（o#3）が指定されたパケット（0__1）が到着すると、パケットバッファ#0は、論理キューVOQ#3にこのパケットを格納する。また、要求振り分け部#0は、論理キューVOQ#3に対応する振り分けポインタの値「1」を参照することにより、この値に対応するスケジューラ部#1にリクエスト通知を送る。リクエスト通知送出後、論理キューVOQ#3に対応する振り分けポインタの値が更新されて「2」になる。

【0056】

T=1において、送出先として出力回線#0（o#0）が指定されたパケット（0__2）が到着すると、パケットバッファ#0は、論理キューVOQ#0にこのパケットを格納する。また、要求振り分け部#0は、論理キューVOQ#0に対応する振り分けポインタの値「0」を参照することにより、この値に対応するスケジューラ部#0にリクエスト通知を送る。リクエスト通知送出後、論理キューVOQ#0に対応する振り分けポインタの値が更新されて「1」になる。

【0057】

T=2において、送出先として出力回線#3（o#3）が指定されたパケット（0__3）が到着すると、パケットバッファ#0は、論理キューVOQ#3にこのパケットを格納する。また、要求振り分け部#0は、T=0において更新された論理キューVOQ#3の振り分けポインタの値「2」を参照することにより、この値に対応するスケジューラ部#2にリクエスト通知を送る。リクエスト通知送出後、論理キューVOQ#3に対応する振り分けポインタの値が更新されて「

3」になる。

【0058】

T=4において、送出先として出力回線#0（o#0）が指定されたパケット（0_4）が到着すると、パケットバッファ#0は、論理キューVOQ#0にこのパケットを格納する。また、要求振り分け部#0は、T=1において更新された論理キューVOQ#0の振り分けポインタの値「1」を参照することにより、この値に対応するスケジューラ部#1にリクエスト通知を送る。リクエスト通知送出後、論理キューVOQ#0に対応する振り分けポインタの値が更新されて「2」になる。

【0059】

このように、各入力バッファ部#0～#3において、一つの論理キューVOQに対応して生成されるリクエスト通知の送出先を α 個のスケジューラ部#0～# $\alpha-1$ に均等に分散させることにより、各入力バッファ部#0～#N-1から送出されるリクエスト通知の送出先をスケジューラ部#0～# $\alpha-1$ のそれぞれに分散させることができる。

【0060】

入力回線毎にリクエスト通知を分散させる場合

図7は、入力回線毎にリクエスト通知を分散させる場合の具体例を示す図である。

各入力バッファ部#0～#3内の要求振り分け部#0～#3は、全ての論理キューVOQに共通の振り分けポインタを有しており、入力回線を介してパケットが到着したときに、この振り分けポインタが示すスケジューラ部#0～#3に対してリクエスト通知を送る。図7では入力バッファ部#0、#3内の丸付き数字によってポインタ値が示されている。

【0061】

以下では、入力回線#0に対応する入力バッファ部#0内の要求振り分け部#0に着目して具体的な動作を説明する。なお、初期状態において、入力回線#0に対応する振り分けポインタの値が「0」に設定されているものとする。

T=0において、送出先として出力回線#3（o#3）が指定されたパケット

(0__1) が到着すると、パケットバッファ # 0 は、論理キュー VOQ # 3 にこのパケットを格納する。また、要求振り分け部 # 0 は、共通の振り分けポインタの値「0」を参照することにより、この値に対応するスケジューラ部 # 0 にリクエスト通知を送る。リクエスト通知送出後、振り分けポインタの値が更新されて「1」になる。

【0062】

T = 1 において、送出先として出力回線 # 0 (o # 0) が指定されたパケット (0__2) が到着すると、パケットバッファ # 0 は、論理キュー VOQ # 0 にこのパケットを格納する。また、要求振り分け部 # 0 は、T = 0 において更新された振り分けポインタの値「1」を参照することにより、この値に対応するスケジューラ部 # 1 にリクエスト通知を送る。リクエスト通知送出後、振り分けポインタの値がさらに更新されて「2」になる。

【0063】

T = 2 以後についても同様であり、パケットが到着していずれかの論理キュー VOQ に格納されると、共通の振り分けポインタの値が参照されてリクエスト通知の送出先となるスケジューラ部 # 0 ~ # 3 が決定される。

このように、各入力バッファ部 # 0 ~ # 3 において、各論理キュー VOQ に対応して生成されるリクエスト通知の送出先を共通の振り分けポインタを用いて α 個のスケジューラ部 # 0 ~ # $\alpha - 1$ に均等に分散させることにより、各入力バッファ部 # 0 ~ # N - 1 から送出されるリクエスト通知の送出先をスケジューラ部 # 0 ~ # $\alpha - 1$ に分散させることができる。

【0064】

単位時間ごとにリクエスト通知を分散させる場合

図 8 は、単位時間毎にリクエスト通知を分散させる場合の具体例を示す図である。

各入力バッファ部 # 0 ~ # 3 内の要求振り分け部 # 0 ~ # 3 は、単位時間 T 毎 (1 パケット時間毎) に更新される振り分けポインタを有しており、入力回線を介してパケットが到着したときに、この振り分けポインタが示すスケジューラ部に対してリクエスト通知を送る。図 8 では入力バッファ部 # 0、# 3 内の丸付き

数字によってポインタ値が示されている。

【0065】

以下では、入力回線 # 0 に対応する入力バッファ部 # 0 内の要求振り分け部 # 0 に着目して具体的な動作を説明する。なお、 $T = 0$ に対応して振り分けポインタの値が「0」に設定されているものとする。

$T = 0$ において、送信先として出力回線 # 3 (o # 3) が指定されたパケット (0_1) が到着すると、パケットバッファ # 0 は、論理キュー VOQ # 3 にこのパケットを格納する。また、要求振り分け部 # 0 は、 $T = 0$ に対応する振り分けポインタの値「0」に基づいてスケジューラ部 # 0 にリクエスト通知を送る。

【0066】

$T = 1$ において、送出先として出力回線 # 0 (o # 0) が指定されたパケット (0_2) が到着すると、パケットバッファ # 0 は、論理キュー VOQ # 0 にこのパケットを格納する。また、要求振り分け部 # 0 は、 $T = 1$ に対応する振り分けポインタの値「1」に基づいてスケジューラ部 3 0 # 1 にリクエスト通知を送る。

【0067】

このようにして時間 T が経過する毎に振り分けポインタの値が巡回的に更新され、その値に対応するスケジューラ部 # 0 ~ # 3 のいずれかにリクエスト通知が送られる。したがって、 $T = 3$ において到着するパケットがない場合には、 $T = 3$ に対応する振り分けポインタの値「3」に基づいてスケジューラ部 # 3 にリクエスト通知が送られずに、次の $T = 4$ において到着するパケット (0_4) に対応するリクエスト通知は、 $T = 4$ に対応する振り分けポインタの値「0」に基づいてスケジューラ部 # 0 に送られる。

【0068】

このように、各入力バッファ部 # 0 ~ # 3 において、各論理キュー VOQ に対応して生成されるリクエスト通知の送出先を、パケットの到着時間 T に対応して更新される振り分けポインタを用いて α 個のスケジューラ部 # 0 ~ # $\alpha - 1$ に均等に分散させることにより、各入力バッファ部 # 0 ~ # $N - 1$ から送出されるリクエスト通知の送出先を各スケジューラ部 # 0 ~ # $\alpha - 1$ に分散させることがで

きる。

【0069】

各スケジューラ部で管理しているリクエスト数を参照してリクエスト通知を分散させる場合

図9は、各スケジューラ部で管理しているリクエスト数を入力バッファ部20が参照してリクエスト通知を分散させる場合の具体例を示す図である。

【0070】

各入力バッファ部#0～#3内の要求振り分け部#0～#3は、パケットが到着したときに、このパケットと同じ入力回線および出力回線（論理キューVOQ）の組合せに対応するリクエスト通知がスケジューラ部#0～#3のそれぞれにいくつあるかを参照し、最もリクエスト数が少ないスケジューラ部をこの到着パケットに対応するリクエスト通知の送出先として決定する。図9に示した4つのスケジューラ部#0～#3の右側には、入力回線#0と論理キューVOQとの組合せに対応するリクエスト数が示されている。以下では、入力バッファ部#0内の要求振り分け部#0に着目して具体的な動作を説明する。

【0071】

$T=0$ において、送信先として出力回線#3（ $\alpha=3$ ）が指定されたパケット（ 0_1 ）が到着すると、パケットバッファ#0は、論理キューVOQ#3にこのパケットを格納する。また、要求振り分け部#0は、4つのスケジューラ部#0～#3内の要求数管理部#0～#3が保持している情報を取得することにより、入力回線#0と送出先の出力回線#3（論理キューVOQ#3）が同じリクエスト数がいくつあるかを調べ、このリクエスト数が最も少ないスケジューラ部#0を選択してリクエスト通知を送る。

【0072】

$T=1$ 以降についても同様であり、対応するリクエスト通知の数がその時点で最も少ないスケジューラ部#0～#3を選択してリクエスト通知を送出する。

このように、各入力バッファ部#0～# $N-1$ において、全てのスケジューラ部#0～# $\alpha-1$ で管理しているリクエスト数取得し、最もリクエスト数が少ないスケジューラ部にリクエスト通知を送ることにより、各入力バッファ部#0

～#N-1から送出されるリクエスト通知を α 個のスケジューラ部#0～# $\alpha-1$ に均等に分散させることができる。

【0073】

リクエスト数の管理機能の分散配備

次に、各スケジューラ部30のみにおいてリクエスト数を管理するのではなく、スケジューラ部30と入力バッファ部20の両方にリクエスト数の管理機能を分散配備した場合のリクエスト振り分け方法について説明する。前提として、各スケジューラ部30内の要求数管理部32は、最低限必要な数だけリクエスト数を管理することができる要求管理機能を有しているものとする。また、各入力バッファ部20内の要求振り分け部24は、論理キューVOQ毎に各スケジューラ部30に送出されたリクエスト数を管理する機能（サブリクエストカウンタによって実現する）と、リクエスト未通知数を管理する機能（リクエストカウンタによって実現する）とを有しているものとする。

【0074】

例えば、各スケジューラ部30内の要求数管理部32が、各論理キューVOQ毎に2個のリクエスト通知を管理することができる機能を有するものとする、入力バッファ部20内の要求振り分け部24からいずれかのスケジューラ部30にリクエスト通知を送出する場合に、対応する論理キューVOQのリクエスト数が2個未満のスケジューラ部30を選択するようにする。全てのスケジューラ部30がこの論理キューVOQについて2個のリクエスト通知を有している場合には、入力バッファ部20は、リクエスト通知の送出は行わずに、この論理キューVOQに対応するリクエストカウンタをインクリメントし、リクエスト数が2未満になったときにリクエスト通知を送出する。

【0075】

図10～図14は、リクエスト数の管理機能を入力バッファ部20とスケジューラ部30に分散配備した場合の具体例を示す図である。これらの図では、説明を簡略化するために、入力回線#0（i#0）に対応する入力バッファ部#0において一つの論理キューVOQ#aに着目した場合の具体例が示されている。なお、スケジューラ部30の数を4とし、スケジューラ部30側における要求管理

数の最大値を2とする。

【0076】

例えば、図10に示すように、 $T = a$ における初期状態として、スケジューラ部#0 (SCH#0) には入力回線#0、論理キューVOQ#aに対応する1個のリクエスト通知が存在し、他のスケジューラ部#1～#3 (SCH#1～SCH#3) にはそれぞれ2個のリクエスト通知が存在するものとする。

【0077】

この状態において、新たにパケットが到着すると、要求振り分け部#0は、リクエスト数が2未満であり、かつ最もリクエスト数が少ないスケジューラ部#0 (SCH#0) にリクエスト通知を送るとともに、スケジューラ部#0に対応するサブリクエストカウンタの値を「2」に更新する。この結果、図11に示すように、 $T = a'$ において、全てのスケジューラ部#0～#3に2個ずつのリクエスト通知が存在することになる。

【0078】

次に、 $T = a + 1$ において新たにパケットが到着しても、この時点で全てのスケジューラ部#0～#3に2個ずつのリクエスト通知があるため、要求振り分け部#0は、新たなリクエスト通知を送出することはせずに、未通知のリクエスト数を管理するために設けられているリクエストカウンタの値をインクリメントする (図12)。

【0079】

その後、スケジューラ部#2から入力バッファ部#0に対してグラント通知が送られてきて、入力バッファ部#0内の論理キュー#aからパケットが読み出されて出力回線#aに送出されると、要求振り分け部#0は、スケジューラ部#2に対応するサブリクエストカウンタの値をデクリメントする。この結果、図13に示すように、 $T = a + 1'$ において、このサブリクエストカウンタの値が「1」に変更される。

【0080】

次に、要求振り分け部#0は、リクエストカウンタの値が1以上であり、しかもスケジューラ部#2に対応するサブリクエストカウンタの値が2未満になった

ことを検出し、このスケジューラ部 # 2 に対してリクエスト通知を送る。また、リクエスト通知後 ($T = a + 1$)、図 1 4 に示すように、要求振り分け部 # 0 は、このスケジューラ部 # 2 に対応するサブリクエストカウンタの値を「2」にインクリメントするとともに、リクエストカウンタの値を「0」にデクリメントする。

【 0 0 8 1 】

このように、各入力バッファ部 # 0 ~ # 3 においてリクエスト数の管理を行うことにより、スケジューラ部 # 0 ~ # 3 のそれぞれでは少ないリクエスト数を管理するだけでスケジューリング処理を行うことができ、スケジューラ部 # 0 ~ # 3 における処理の負担を軽減することができる。

【 0 0 8 2 】

スケジューラ部の冗長構成

次に、必要数以上の数のスケジューラ部 3 0 を備えた場合の冗長構成について説明する。

図 1 5 は、冗長系のスケジューラ部 3 0 を非固定にした場合の具体例を示す図である。例えば、各周期が 4 パケット時間に対応しており、スケジューラ部 3 0 の数を「5」とした場合の具体例が示されている。

【 0 0 8 3 】

5 つのスケジューラ部 # 0 ~ # 4 (SCH # 0 ~ SCH # 4) が正常に動作している場合には、それぞれにおいて独立にスケジューリング処理が行われ、全てのスケジューラ部 # 0 ~ # 4 のスケジューリング結果が巡回的に使用される。スケジューラ部 # 0 ~ # 4 のそれぞれは、4 パケット時間で 1 回のスケジューリング処理を行うことが可能であるが、図 1 5 に示した例では、5 パケット時間毎に各スケジューラ部 # 0 ~ # 4 のスケジューリング結果が使用されることになる。

【 0 0 8 4 】

その後、周期 C においてスケジューラ部 # 1 (SCH # 1) に障害が発生すると、各入力バッファ部 2 0 では、次の周期 D からはこの障害が発生したスケジューラ部 # 1 によるスケジューリング結果は使用せずに、他の 4 つのスケジューラ部 # 0、# 2 ~ # 4 のスケジューリング結果のみを巡回的に使用する。

【 0 0 8 5 】

また、周期Eにおいてスケジューラ部 # 1 の障害が復旧した場合には、各入力バッファ部 2 0 は、次の周期Fからは再び全てのスケジューラ部 # 0 ~ # 4 のスケジューリング結果を巡回的に使用する。

図 1 6 は、冗長系のスケジューラ部 3 0 を固定にした場合の動作例を示す図である。例えば、スケジューラ部 # 4 (SCH # 4) は、冗長系 (予備系) として備わったものであり、障害等が発生した場合にのみ使用される。

【 0 0 8 6 】

4 つのスケジューラ部 # 0 ~ # 3 が正常に動作している場合には、それぞれにおいて独立にスケジューリング処理が行われ、4 つのスケジューラ部 # 0 ~ # 3 のそれぞれのスケジューリング結果が巡回的に利移用される。この場合の動作は、冗長系のスケジューラ部 # 4 を有しない場合の動作と同じである。

【 0 0 8 7 】

次に、周期Cにおいてスケジューラ部 # 1 (SCH # 1) に障害が発生すると、このスケジューラ部 # 1 の代わりに冗長系のスケジューラ部 # 4 が使用される。同じ機能を有するスケジューラ部 # 1 とスケジューラ部 # 4 とを入れ替えただけであるため、障害発生前と同様のスケジューリング処理を継続することができる。

【 0 0 8 8 】

また、周期Eにおいてスケジューラ部 # 1 の障害が復旧した場合には、再び冗長系のスケジューラ部 # 4 が待機状態になり、スケジューラ部 # 1 が使用されてスケジューリング処理が継続される。

このように、1 つのスケジューラ部 3 0 がスケジューリング処理に要するパケット時間数よりも多くのスケジューラ部 3 0 を備えた冗長構成を採用することにより、障害が発生した場合であっても、この冗長分 (図 1 5、図 1 6 に示す例では 1 個) のスケジューラ部 3 0 に代わりの処理を行わせることが可能になるため、遅延を来すことなくスケジューリング処理を継続することができる。

【 0 0 8 9 】

スケジューラ部の拡張性

次に、本実施形態のパケットスイッチ 1 0 0 に収容される入力回線の数や出力回線の数に応じてスケジューラ部 3 0 の数を変更する場合の拡張性について説明する。なお、以下では説明を簡略化するために、最大のスイッチ規模を 4×4 （入力回線数 \times 出力回線数）とし、そのときのスケジューラ部 3 0 の数を 4、各スケジューラ部 3 0 がスケジューリング処理に要する時間を 4 パケット時間（ $= 4 T_p$ ）とする。また、このようなスケジューラ部 3 0 を用いてスイッチ規模が 2×2 のパケットスイッチのスケジューリング処理を行う場合を考えるものとする。

【 0 0 9 0 】

図 1 7 は、最大規模 4×4 のパケットスイッチを構成した場合のスケジューリング処理の概要を示す図である。例えば、4 個のスケジューラ部の中の一つに着目し、ラウンドロビンを用いてスケジューリング処理を行う場合の具体例が示されている。

【 0 0 9 1 】

最初の周期 # 0 において、入力回線 # 0、# 1、# 2、# 3 の順にスケジューリング処理が行われる。スケジューラ部 # 0 ~ # 3 のそれぞれがスケジューリング処理に必要な時間は 4 パケット時間（ 1τ ）であるため、4 パケット時間毎に各スケジューラ部によって 1 つのスケジューリング結果が得られる。図 1 7 に示した例では、周期 # 0 において、以下に示すスケジューリング結果が得られる。

【 0 0 9 2 】

- ・ 入力回線 # 0 (i # 0) \rightarrow 出力回線 # 0 (o # 0)
- ・ 入力回線 # 1 (i # 1) \rightarrow 出力回線 # 1 (o # 1)
- ・ 入力回線 # 2 (i # 2) \rightarrow 出力回線 # 2 (o # 2)
- ・ 入力回線 # 3 (i # 3) \rightarrow 出力回線 # 3 (o # 3)

また、次の周期 # 1 において、入力回線 # 1、# 2、# 3、# 0 の順にスケジューリング処理が行われ、以下に示すスケジューリング結果が得られる。

【 0 0 9 3 】

- ・ 入力回線 # 1 (i # 1) \rightarrow 出力回線 # 1 (o # 1)
- ・ 入力回線 # 2 (i # 2) \rightarrow 出力回線 # 3 (o # 3)

- ・ 入力回線 # 3 (i # 3) → 出力回線 # 2 (o # 2)
- ・ 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

周期 # 2 以降についても同様であり、スケジューリング処理の対象となる入力回線の処理順番を巡回させながら、各入力回線に対応する出力回線が各周期毎に決定される。

【 0 0 9 4 】

このように、パケットスイッチ全体では、4つのスケジューラ部 # 0 ~ # 3 によって4パケット時間毎に4つのスケジューリング結果が得られる。平均すると、1パケット時間毎に1つのスケジューリング結果が得られることになり、スループットの劣化がないことがわかる。

【 0 0 9 5 】

図 1 8 は、2 × 2 の小規模のパケットスイッチを構成した場合のスケジューリング処理の概要を示す図である。図 1 8 には、図 1 7 に示した最大構成時と同じスケジューラ部を2つ用いた場合の具体例が示されている。

最初の周期 # 0 において、入力回線 # 0、# 1 の順にスケジューリング処理が行われる。2つのスケジューラ部 # 0、# 1 のそれぞれは、4本の入力回線のそれぞれに対応する出力回線を決定する処理能力を有するが、図 1 8 に示す例では、処理対象となる入力回線は2本であるため、4パケット時間 (1 τ) 毎にこれら2本の入力回線に対するスケジューリング処理が実施され、以下に示すスケジューリング結果が得られる。

【 0 0 9 6 】

- ・ 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)
- ・ 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

また、周期 # 1 において、入力回線 # 1、# 0 の順にスケジューリング処理が行われ、以下に示すスケジューリング結果が得られる。

【 0 0 9 7 】

- ・ 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)
- ・ 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

周期 # 2 以降についても同様であり、スケジューリング処理の対象となる入力

回線の処理順番を巡回させながら、各入力回線に対応する出力回線が各周期毎に決定される。

【 0 0 9 8 】

このように、入力回線数や出力回線数を減少させたことに伴ってスケジューラ部の数を減らしただけではスループットの低下を招くことになる。すなわち、2つのスケジューラ部 # 0、# 1 のそれぞれにおいてスケジューリング処理を行うために必要な時間は4 パケット時間であることに変わりはなく、パケットスイッチ # 0、# 1 の全体では、4 パケット時間毎に2つのスケジューリング結果が得られることになる。したがって、4 パケット時間の内の半分については、スケジューラ部 # 0、# 1 から入力バッファ部 # 0、# 1 のそれぞれに対してパケットの送出指示を行うことができず、スループットが半分に低下する。

【 0 0 9 9 】

図 1 9 は、スループットの低下を防止した小規模のパケットスイッチにおけるスケジューリング処理の概要を示す図である。2 × 2 のパケットスイッチにおいて2個のスケジューラ部によってスループットを低下させることなくスケジューリング処理を行うためには、図 1 9 に示すように、スケジューラ部 # 0、# 1 のそれぞれにおけるスケジューリング処理の所要時間を2 パケット時間に設定する。処理対象となる入力回線数や出力回線数のそれぞれが半分になっているため、この所要時間の変更は可能である。

【 0 1 0 0 】

最初の周期 # 0 の前半において、入力回線 # 0、# 1 の順にスケジューリング処理が行われ、以下に示すスケジューリング結果が得られる。

- ・ 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)
- ・ 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

また、この周期 # 0 の後半において、入力回線 # 1、# 0 の順にスケジューリング処理が行われ、以下に示すスケジューリング結果が得られる。

【 0 1 0 1 】

- ・ 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)
- ・ 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

周期 # 1 以降についても同様であり、前半と後半のそれぞれにおいて別々にスケジューリング処理が行われる。これにより、パケットスイッチ全体では、2 パケット時間毎に2つのスケジューリング結果が得られる。平均すると、1 パケット時間毎に1つのスケジューリング結果が得られることになり、スループットの劣化がないことがわかる。

【0 1 0 2】

図 2 0 は、スループットの低下を防止した小規模のパケットスイッチにおけるスケジューリング処理の概要を示す図である。例えば、2 個のスケジューラ部のそれぞれのスケジューリング時間は4 パケット時間のままで、スループットを低下させることなくスケジューリング処理を行う場合の動作手順が示されている。

【0 1 0 3】

具体的には、スケジューリング時間を変えない代わりに、各周期の中で複数のスケジューリング結果を導く方法が採用されている。すなわち、 2×2 の小規模スイッチの場合には、入力回線 # 2、# 3 と出力回線 # 2、# 3 に対するスケジューリング処理は行われないので、これらの未使用回線分の処理時間を利用して、1 回のスケジューリング周期（1 τ 、4 パケット時間）の中で2 個のスケジューリング結果を導き出す。

【0 1 0 4】

図 2 0 において、下線が付された入力回線 # 0、# 1 (i # 0、i # 1) は、それぞれ入力回線 # 2、# 3 の代わりにスケジューリング処理が行われることを示している。また、出力回線 # 0、# 1 (o # 0、o # 1) は、スケジューリング結果としてパケットの初出先として出力回線 # 2、# 3 が得られた場合に、それぞれを入力回線 # 0、# 1 に置き換えたものであることを示している。

【0 1 0 5】

最初の周期 # 0 において、入力回線 # 0、# 1、# 2、# 3 の順にスケジューリング処理が行われる。スケジューラ部 # 0、# 1 のそれぞれがスケジューリング処理に必要な時間は4 パケット時間（1 τ ）であるため、4 パケット時間毎に各スケジューラ部によって4 本の入力回線 # 0 ~ # 3 に対応する1つのスケジューリング結果が得られる。ところで、上述したように、入力回線 # 2、# 3 は実

際に存在せず、これらの代わりに入力回線 # 0、# 1 についてスケジューリング処理が行われるため、結果的に、スケジューラ部 # 0、# 1 のそれぞれは、4 パケット時間毎に 2 つのスケジューリング結果を導き出すことになる。図 2 0 に示した例では、周期 # 0 において、以下に示すスケジューリング結果が得られる。

【0 1 0 6】

結果 1 : 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

結果 1 : 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

結果 2 : 入力回線 # 0 (i # 0) → 出力回線 # 0 (o # 0)

結果 2 : 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

また、次の周期 # 1 において、入力回線 # 1、# 2、# 3、# 0 の順にスケジューリング処理が行われ、以下に示すスケジューリング結果が得られる。

【0 1 0 7】

結果 1 : 入力回線 # 1 (i # 1) → 出力回線 # 1 (o # 1)

結果 2 : 入力回線 # 0 (i # 0) → 出力回線 # 1 (o # 1)

結果 2 : 入力回線 # 1 (i # 1) → 出力回線 # 0 (o # 0)

結果 1 : 入力回線 # 0 (i # 0) → 出力回線 # 1 (o # 1)

このように、2 つのスケジューラ部 # 0、# 1 のそれぞれは、4 パケット時間の 1 周期の中で 2 つのスケジューリング結果を導き出しており、パケットスイッチ全体では、4 パケット時間毎に 2 つのスケジューラ部 # 0、# 1 によって合計で 4 つのスケジューリング結果を得ることができる。平均すると、1 パケット時間毎に 1 つのスケジューリング結果が得られることになり、スループットの劣化がないことがわかる。

【0 1 0 8】

このように、本実施形態のパケットスイッチ 1 0 0 では、入力回線数や出力回線数に応じたスケジューラ部の数を設定するとともに、各スケジューラ部におけるスケジューリング処理の所要時間を変更することにより、スループットの低下を招くことなく、最適なスケジューリング処理を行うことができる。したがって、無駄な構成を低減するとともに拡張性に優れたパケットスイッチを実現することができる。

【0 1 0 9】

次に、各スケジューラ部によるスケジューリング所要時間に対応する1周期の中で、複数のスケジューリング結果を導き出すスケジューリング処理の具体例を説明する。

図21は、1周期の中で複数のスケジューリング結果を導き出すスケジューリング処理の具体例を示す図であり、図20に示した周期#0におけるスケジューリング処理の内容が模式的に示されている。スケジューラ部#0、#1のそれぞれは、入力回線と論理キューVOQの組み合わせに対応するリクエスト通知の有無を管理している。例えば、入力回線#0については、出力回線#0と出力回線#1の両方に対するリクエスト通知が存在していることがわかる。この例では、2×2の小規模スイッチを想定しているため、入力回線#2、#3および出力回線#2、#3は使用されていない。このため、入力回線#0、#1のいずれかから出力回線#2、#3のいずれかにパケットを送出する旨の内容のリクエスト通知は存在しない。

【0 1 1 0】

各スケジューラ部は、1周期の中で複数のスケジューリング結果を導き出すために、未使用回線に対応するスケジューリング処理を実施する。具体的には、入力回線#2に対応する処理を行うブロックに入力回線#0のリクエスト通知の内容をセットする。このとき、入力回線#0と出力回線#0とを組み合わせたリクエスト通知に対しては、入力回線#2と出力回線#2の組み合わせに変更した情報が代わりにセットされる。入力回線#0と出力回線#1とを組み合わせたリクエスト通知に対しては、入力回線#2と出力回線#3の組み合わせに変更した情報が代わりにセットされる。以下に、組合せを変更する場合の具体的な関係を示す。

【0 1 1 1】

i # 2 - o # 0	要求有無情報	←	i # 0 - o # 2	要求有無情報
i # 2 - o # 1	要求有無情報	←	i # 0 - o # 3	要求有無情報
i # 2 - o # 2	要求有無情報	←	i # 0 - o # 0	要求有無情報
i # 2 - o # 3	要求有無情報	←	i # 0 - o # 1	要求有無情報

i # 3 - o # 0 要求有無情報 ← i # 1 - o # 2 要求有無情報

i # 3 - o # 1 要求有無情報 ← i # 1 - o # 3 要求有無情報

i # 3 - o # 2 要求有無情報 ← i # 1 - o # 0 要求有無情報

i # 3 - o # 3 要求有無情報 ← i # 1 - o # 1 要求有無情報

ここで、i # 2等は入力回線番号を、o # 0等は出力回線番号を、要求有無情報はリクエスト要求の有無を示す情報をそれぞれ示している。

【0 1 1 2】

次に、図 2 1 を参照しながら具体的なスケジューリング処理の手順を説明する。

第 1 のスケジューリング対象は入力回線 # 0 であり、これに対するスケジューリング処理は、未確定でしかもリクエスト通知が存在する出力回線の中から一つを選択することにより行われる。図 2 1 に示した例では、出力回線 # 0 が選択される。このようにして選択された出力回線 # 0 は、確定済みとなって未確定情報に反映され、次のスケジューリング対象に対するスケジューリング処理に受け渡される。

【0 1 1 3】

第 2 のスケジューリング対象は入力回線 # 1 であり、これに対するスケジューリング処理も、未確定でしかもリクエスト通知が存在する出力回線の中から一つを選択することにより行われる。図 2 1 に示した例では、既に出力回線 # 0 が確定済みとなっており、出力回線 # 1 が選択される。このようにして選択された出力回線 # 1 は、確定済みとなって未確定情報に反映され、次のスケジューリング対象に対するスケジューリング処理に受け渡される。

【0 1 1 4】

第 3 のスケジューリング対象は入力回線 # 2 である。上述したように、この入力回線 # 2 に対応して設定されている要求有無情報は、入力回線 # 0 に対応するものであるため、実際には入力回線 # 0 についてスケジューリング処理を行うことになる。入力回線 # 2 に対するスケジューリング処理も、未確定でしかもリクエスト通知が存在する出力回線の中から一つを選択することにより行われる。図 2 1 に示した例では、この段階で、未確定でしかもリクエスト通知が存在する出

力回線は、出力回線 # 2 と出力回線 # 3 であり、ここでは出力回線 # 2 が選択される。このようにして選択された出力回線 # 2 は、確定済みとなって未確定情報に反映され、次のスケジューリング対象に対するスケジューリング処理に反映される。

【0115】

なお、出力回線 # 2 は実際には未使用であり、このような未使用回線に対するスケジューリング処理においては、その結果選択された出力回線番号をそのまま入力バッファ部 20 側に通知するのではなく、読み替え処理を行った後の出力回線番号が通知される。すなわち、入力回線 # 2 は入力回線 # 0 に、出力回線 # 2 は出力回線 # 0 に読み替えて通知が行われる。以下に、読み替え処理が行われる入力回線と出力回線の関係を示す。

【0116】

i # 2 - o # 2 → i # 0 - o # 0

i # 2 - o # 3 → i # 0 - o # 1

i # 3 - o # 2 → i # 1 - o # 0

i # 3 - o # 3 → i # 1 - o # 1

同様にして、第 4 のスケジューリング対象である入力回線 # 3 についてもスケジューリング処理が行われ、その結果出力回線 # 3 が選択される。入力回線 # 3、出力回線 # 3 はともに未使用回線であるため、実際には、それぞれが入力回線 # 1、出力回線 # 1 に読み替えられて入力バッファ部 20 に通知される。

【0117】

このように、未使用回線を利用してスケジューリング処理を行うことにより、スケジューリング処理の 1 周期の中で複数の結果を得ることができる。

(付記 1) N 本の入力回線のそれぞれに対応して設けられ、対応する前記入力回線を介して入力される packets を格納する N 個の入力バッファ部と、

それぞれが独立したスケジューリング処理を行うことにより、前記 N 個の入力バッファ部のそれぞれに格納された前記 packets の送出先となる M 本の出力回線のいずれかを決定する α 個のスケジューラ部と、

前記 N 個の入力バッファ部のそれぞれから出力される前記 packets を、前記ス

ケジューラ部によって決定された送出先となる前記出力回線に出力するスイッチ部と、

を備え、前記N個の入力バッファ部において前記 α 個のスケジューラ部によるスケジューリング処理の結果を巡回的に使用することを特徴とするパケットスイッチ。

(付記2) 付記1において、

前記スケジューラ部によるスケジューリング処理は、前記N個の入力バッファ部から送られてくるスケジューリング要求通知に対応して行われており、

前記N個の入力バッファ部のそれぞれの、前記スケジューリング要求通知の送出先となる前記スケジューラ部を分散させることを特徴とするパケットスイッチ。

(付記3) 付記2において、

前記入力バッファ部は、前記M本の出力回線のそれぞれを送出先とするパケットを格納するM個のキューを有しており、これらM個のキューのそれぞれ毎に前記スケジューリング要求通知の送出先となる前記スケジューラ部を巡回させることを特徴とするパケットスイッチ。

(付記4) 付記2において、

前記入力バッファ部は、前記入力回線毎に前記スケジューリング要求通知の送出先となる前記スケジューラ部を巡回させることを特徴とするパケットスイッチ。

(付記5) 付記2において、

前記入力バッファ部は、単位時間毎に前記スケジューリング要求通知の送出先となる前記スケジューラ部を巡回させることを特徴とするパケットスイッチ。

(付記6) 付記2において、

前記入力バッファ部は、前記 α 個のスケジューラ部のそれぞれについて未処理の前記スケジューリング要求通知の数を調べ、この数が少ない前記スケジューラ部に次の前記スケジューリング要求通知を送ることを特徴とするパケットスイッチ。

(付記7) 付記2において、

前記入力バッファ部は、前記 α 個のスケジューラ部のそれぞれに送出した前記スケジューリング要求通知の数を管理しており、この数が所定数に達した前記スケジューラ部に対する前記スケジューリング要求通知の送出動作を、この数が前記所定数より少なくなるまで遅らせることを特徴とするパケットスイッチ。

(付記 8) 付記 1 ～ 7 のいずれかにおいて、

前記スケジューラ部が前記スケジューリング処理に要する時間が、前記パケットの最小送出間隔の L 倍であるときに、

前記スケジューラ部の数 α は、前記倍数 L 以上の値に設定されていることを特徴とするパケットスイッチ。

(付記 9) 付記 8 において、

$L - \alpha$ が 1 以上に設定されており、

前記 N 個の入力バッファ部において前記 α 個のスケジューラ部の全てのスケジューリング処理の結果を巡回的に使用することを特徴とするパケットスイッチ。

(付記 1 0) 付記 8 において、

$L - \alpha$ が 1 以上に設定されており、

$L - \alpha$ 個の前記スケジューラ部を冗長系として用いるとともに、この冗長系以外の前記スケジューラ部に障害が発生したときに、代わりに冗長系の前記スケジューラ部を用いることを特徴とするパケットスイッチ。

(付記 1 1) 付記 1 ～ 1 0 のいずれかにおいて、

前記入力回線の数 N と前記出力回線の数 M に応じて、前記スケジューラ部の数 α および前記スケジューリング処理の時間を可変に設定することを特徴とするパケットスイッチ。

(付記 1 2) 付記 1 ～ 1 0 のいずれかにおいて、

前記スケジューラ部は、未使用回線を含む前記スケジューリング処理を行っており、実際に使用している前記入力回線および前記出力回線と前記未使用回線との間の読み替え処理を行うことにより、1 回の前記スケジューリング処理によって複数のスケジューリング処理結果を得ることを特徴とするパケットスイッチ。

【 0 1 1 8 】

【発明の効果】

上述したように、本発明によれば、複数のスケジューラ部は互いに独立してスケジューリング処理を行っているため、これらの間の競合制御等を行う必要がなく、このための信号線による信号の遅延等の問題もないことから、実装時の制約が大幅に緩和される。また、スケジューラ部は必要な数だけ用意すればよいため無駄な構成がなく、しかも後に追加することも可能であるため、拡張性に富んだパケットスイッチを実現することができる。

【図面の簡単な説明】

【図 1】

一実施形態のパケットスイッチの構成図である。

【図 2】

入力バッファ部の詳細構成を示す図である。

【図 3】

スケジューラ部の詳細構成を示す図である。

【図 4】

負荷分散スケジューリング処理の具体例を示す図である。

【図 5】

負荷分散スケジューリング処理の他の具体例を示す図である。

【図 6】

論理キュー毎にリクエスト通知を分散させる場合の具体例を示す図である。

【図 7】

入力回線毎にリクエスト通知を分散させる場合の具体例を示す図である。

【図 8】

単位時間毎にリクエスト通知を分散させる場合の具体例を示す図である。

【図 9】

各スケジューラ部で管理しているリクエスト数を入力バッファ部が参照してリクエスト通知を分散させる場合の具体例を示す図である。

【図 1 0】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図である。

【図 1 1】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図である。

【図 1 2】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図である。

【図 1 3】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図である。

【図 1 4】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図である。

【図 1 5】

冗長系のスケジューラ部を非固定にした場合の具体例を示す図である。

【図 1 6】

冗長系のスケジューラ部を固定にした場合の具体例を示す図である。

【図 1 7】

最大規模 4×4 のパケットスイッチを構成した場合のスケジューリング処理の概要を示す図である。

【図 1 8】

2×2 の小規模のパケットスイッチを構成した場合のスケジューリング処理の概要を示す図である。

【図 1 9】

スループットの低下を防止した小規模のパケットスイッチにおけるスケジューリング処理の概要を示す図である。

【図 2 0】

スループットの低下を防止した小規模のパケットスイッチにおけるスケジューリング処理の概要を示す図である。

【図 2 1】



1 周期の中で複数のスケジューリング結果を導き出すスケジューリング処理の具体例を示す図である。

【図 2 2】

スケジューリング機能を分散して配備した従来のパケットスイッチの構成図である。

【図 2 3】

スケジューリング機能を集中して配備した従来のパケットスイッチの構成図である。

【符号の説明】

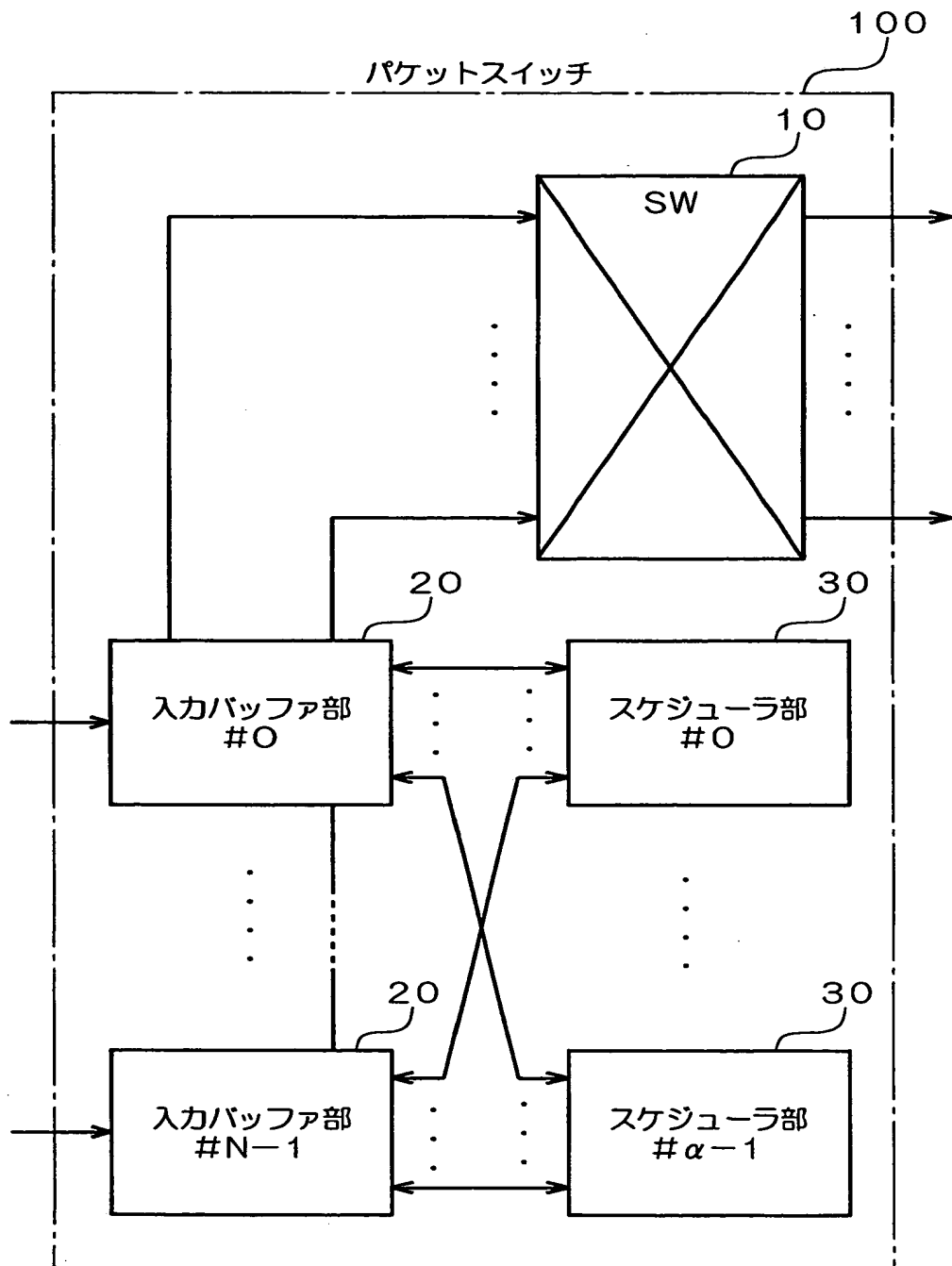
- 1 0 スイッチ部 (SW)
- 2 0 入力バッファ部
- 2 2 パケットバッファ
- 2 4 要求振り分け部
- 2 6 読み出し指示部
- 3 0 スケジューラ部
- 3 2 要求数管理部
- 3 4 スケジューリング制御部
- 1 0 0 パケットスイッチ



【書類名】 図面

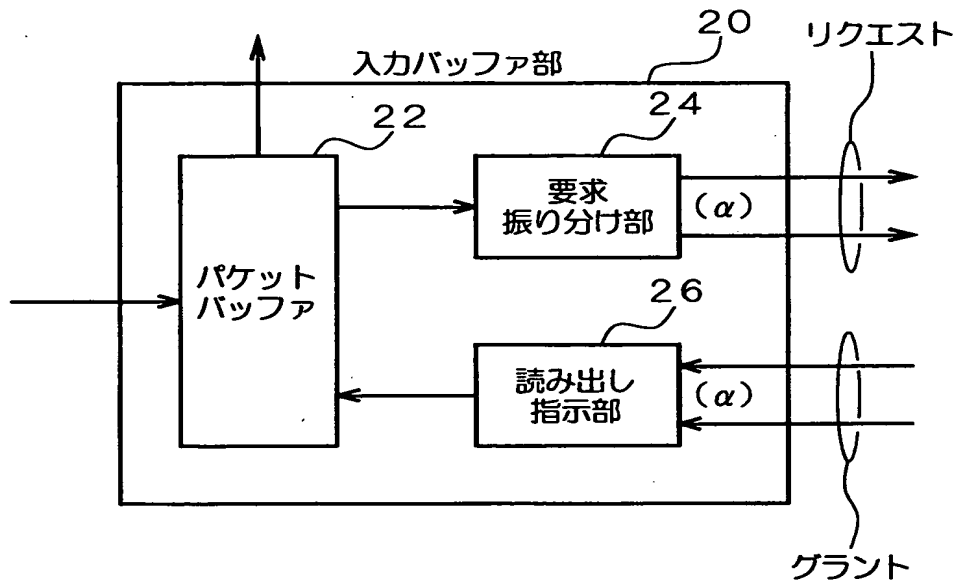
【図1】

一実施形態のパケットスイッチの構成図



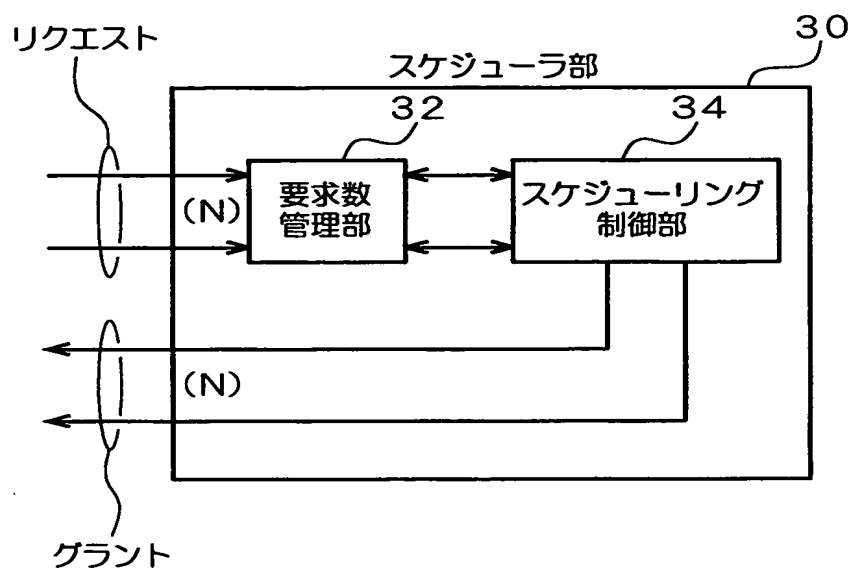
【図 2】

入力バッファ部の構成図



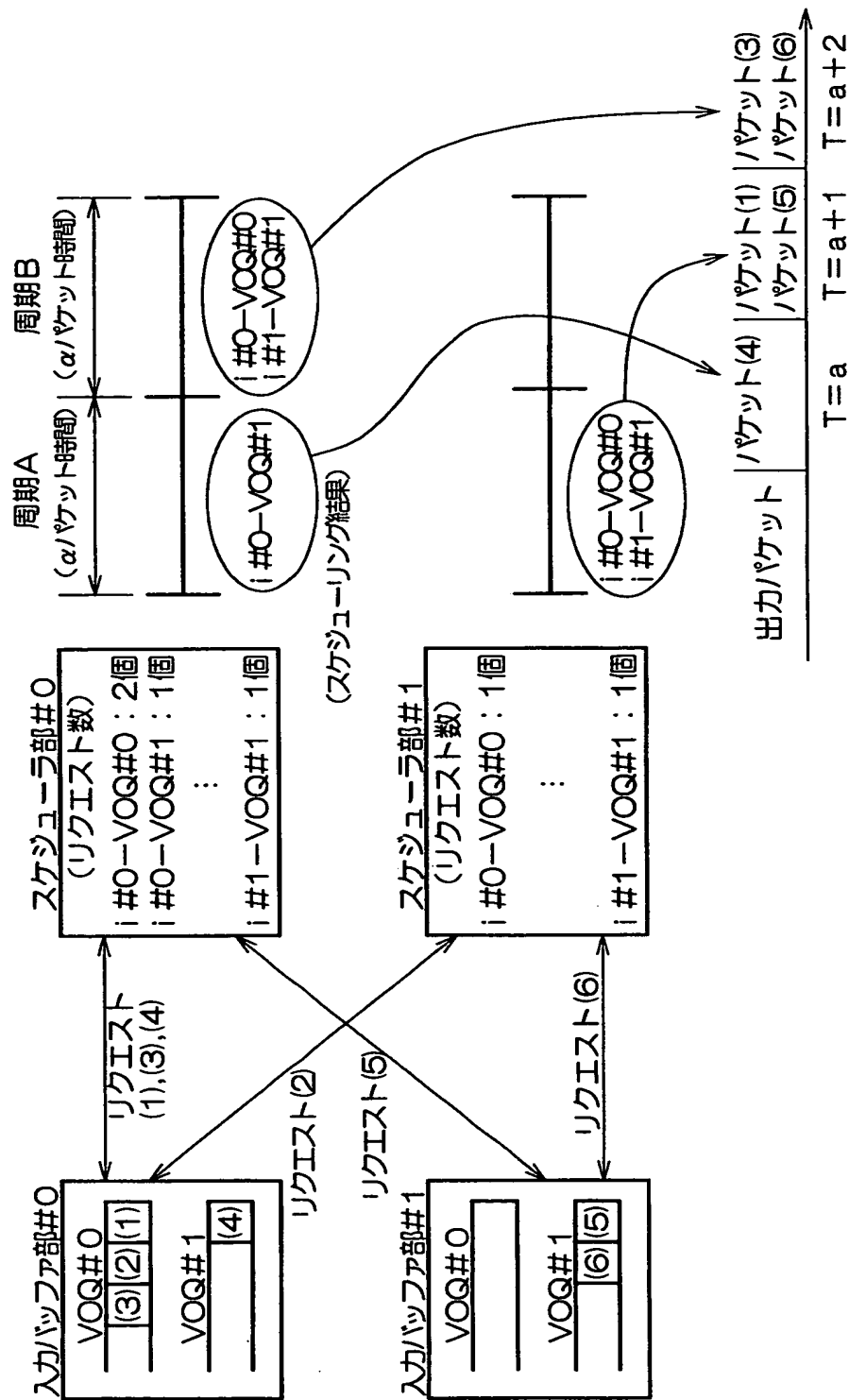
【図 3】

スケジューラ部の構成図



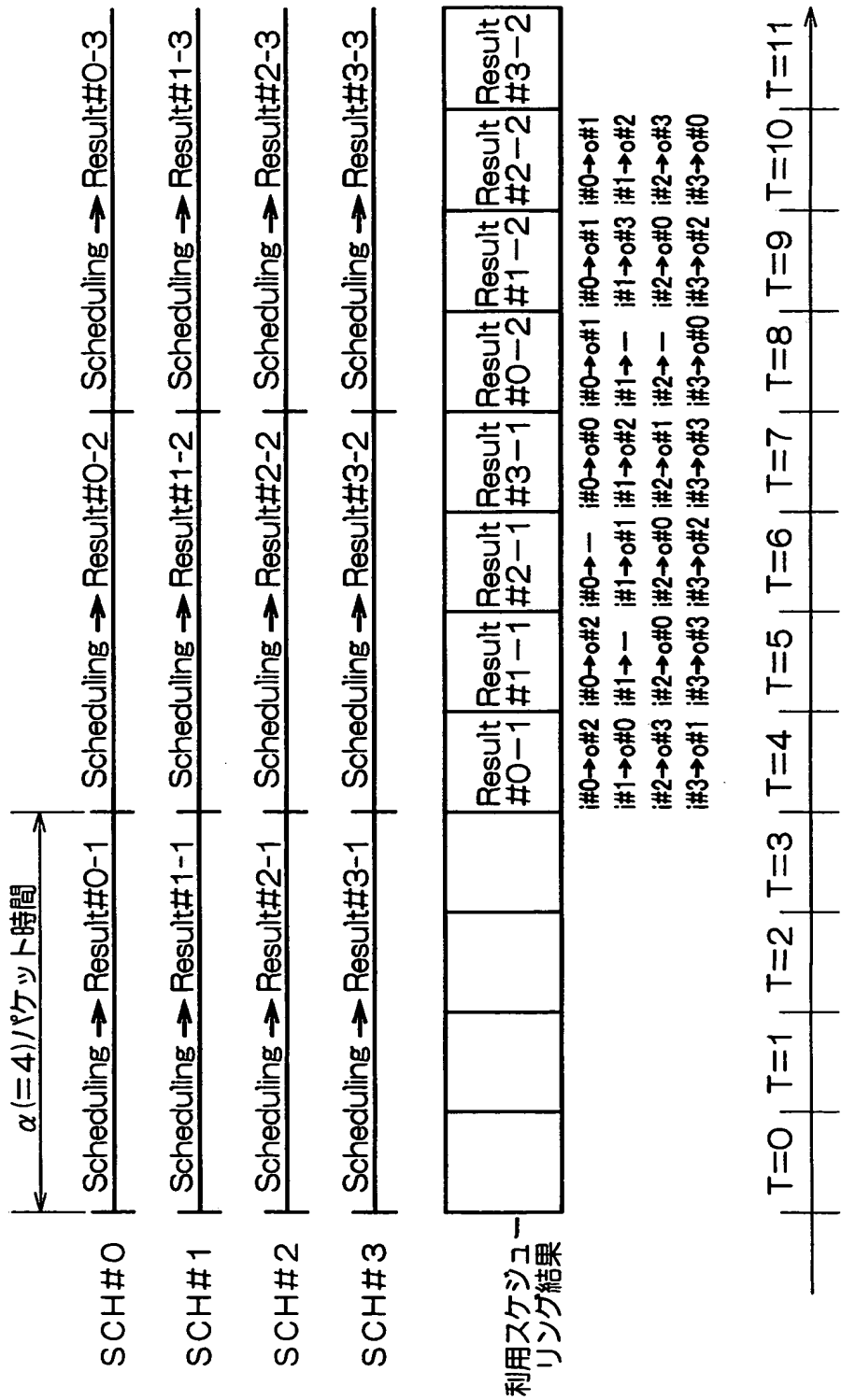
【図 4】

負荷分散スケジューリング処理の具体例を示す図



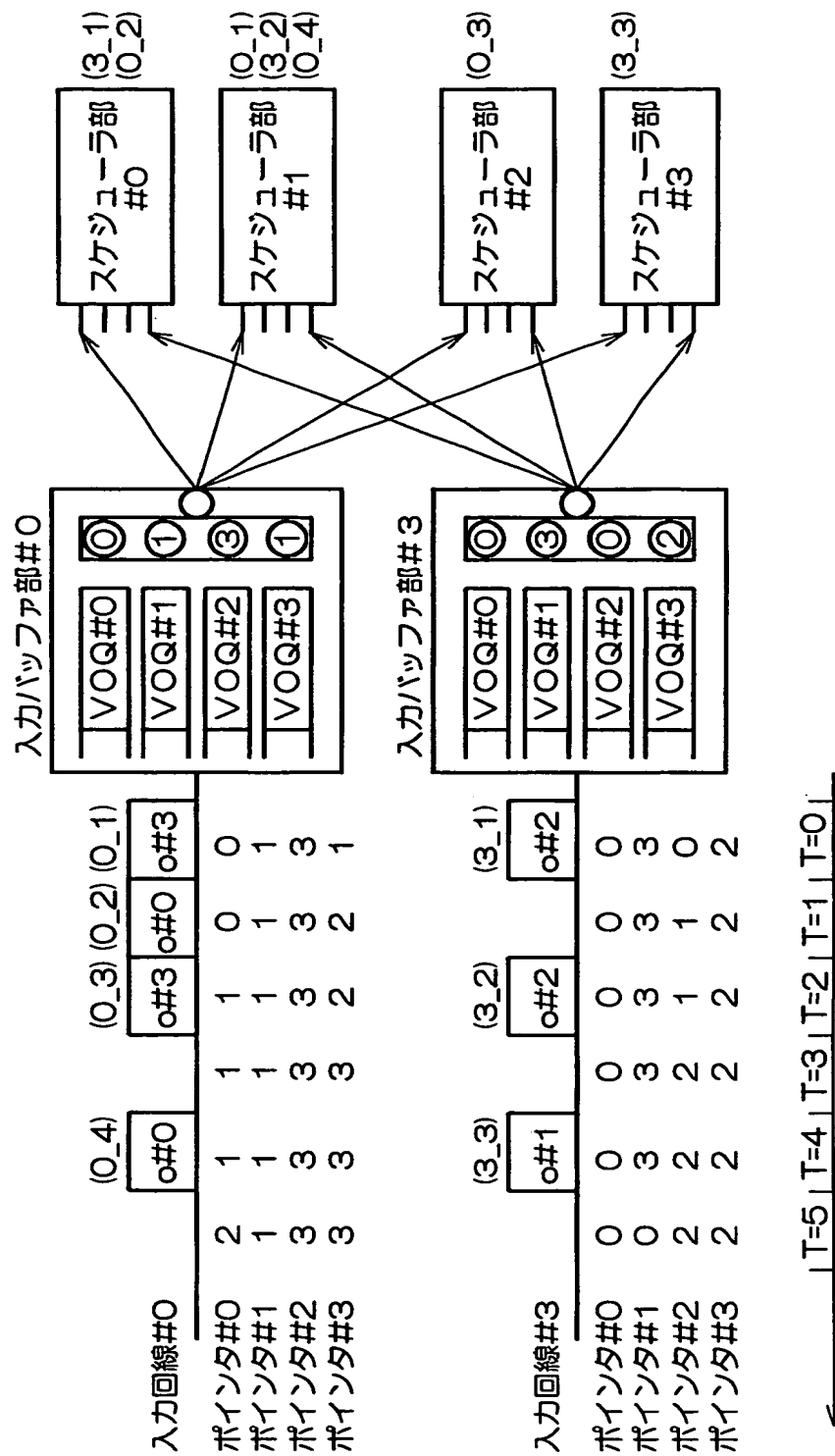
【図 5】

負荷分散スケジューリング処理の具体例を示す図



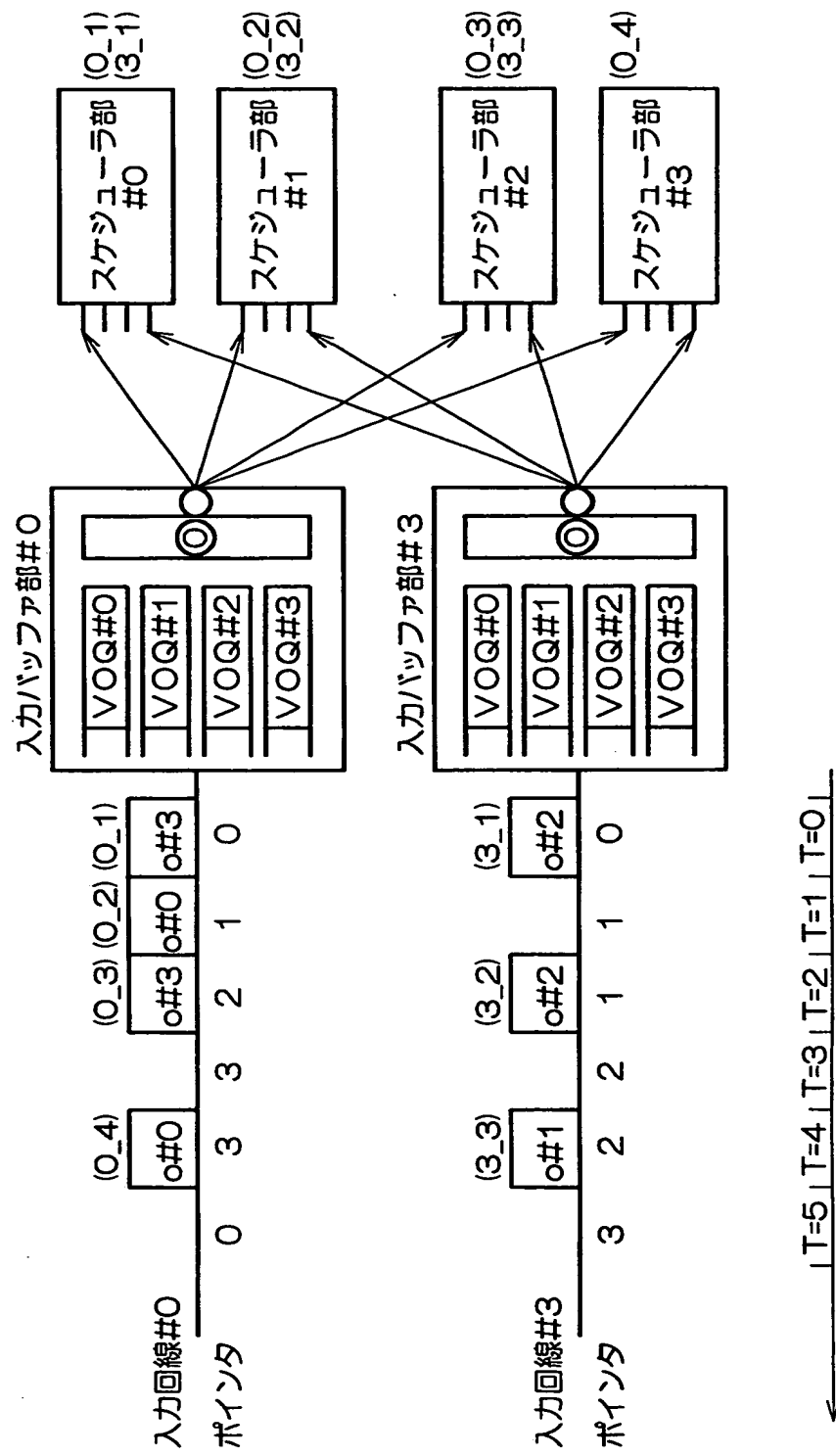
【図 6】

論理キュー毎にリクエスト通知を分散させる場合の具体例を示す図



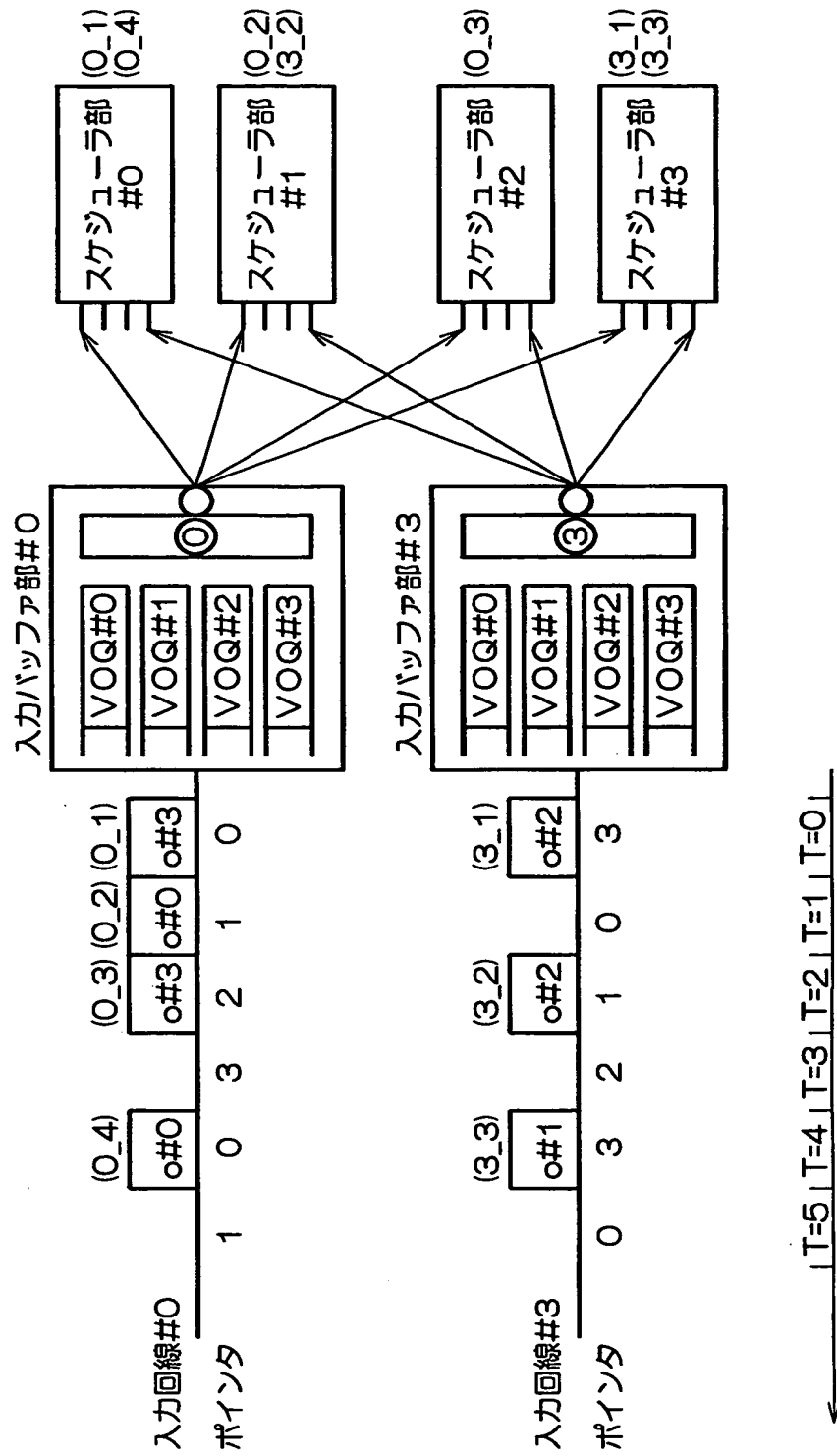
【図 7】

入力回線毎にリクエスト通知を分散させる場合の具体例を示す図

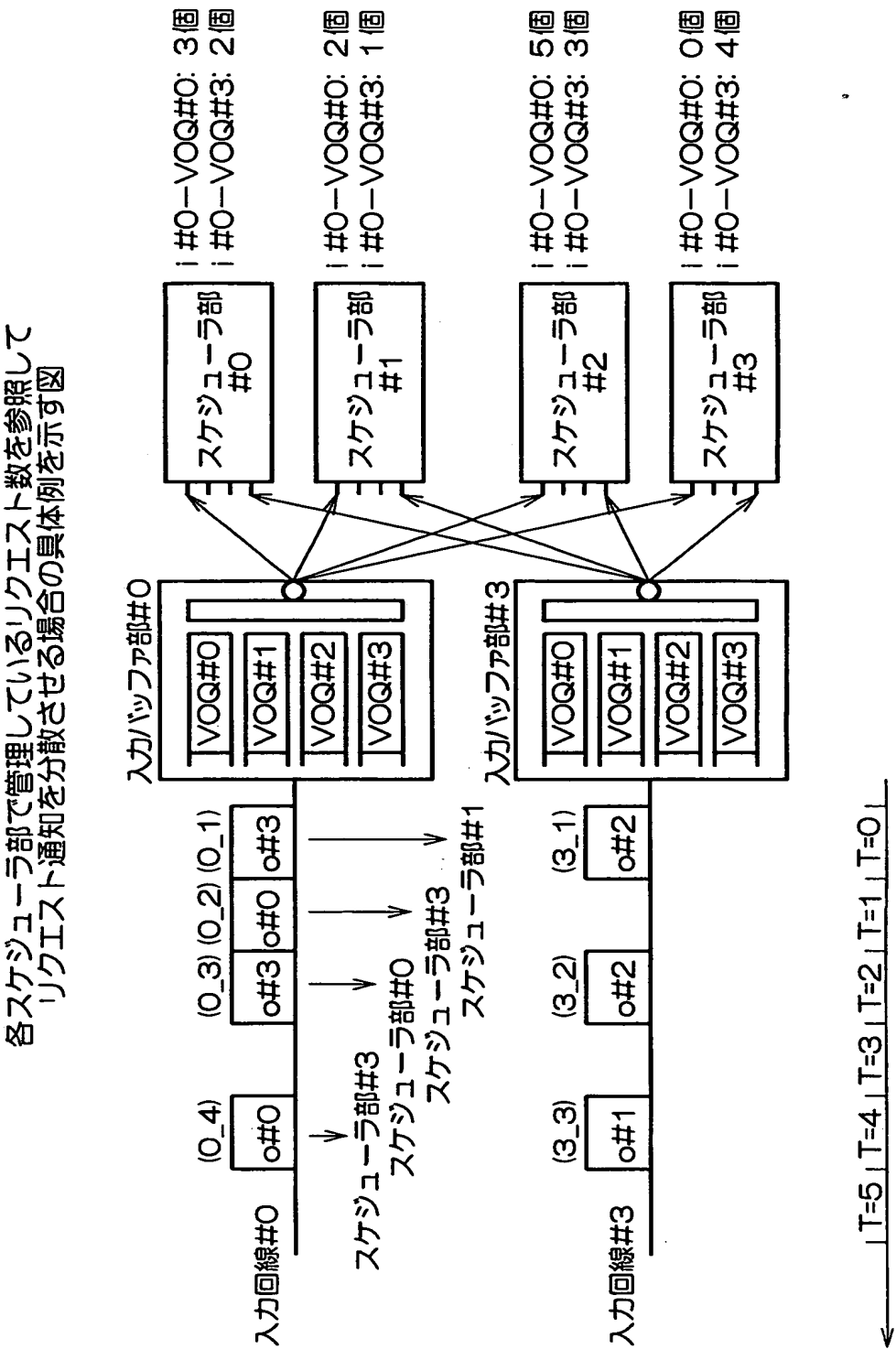


【図 8】

単位時間毎にリクエスト通知を分散させる場合の具体例を示す図

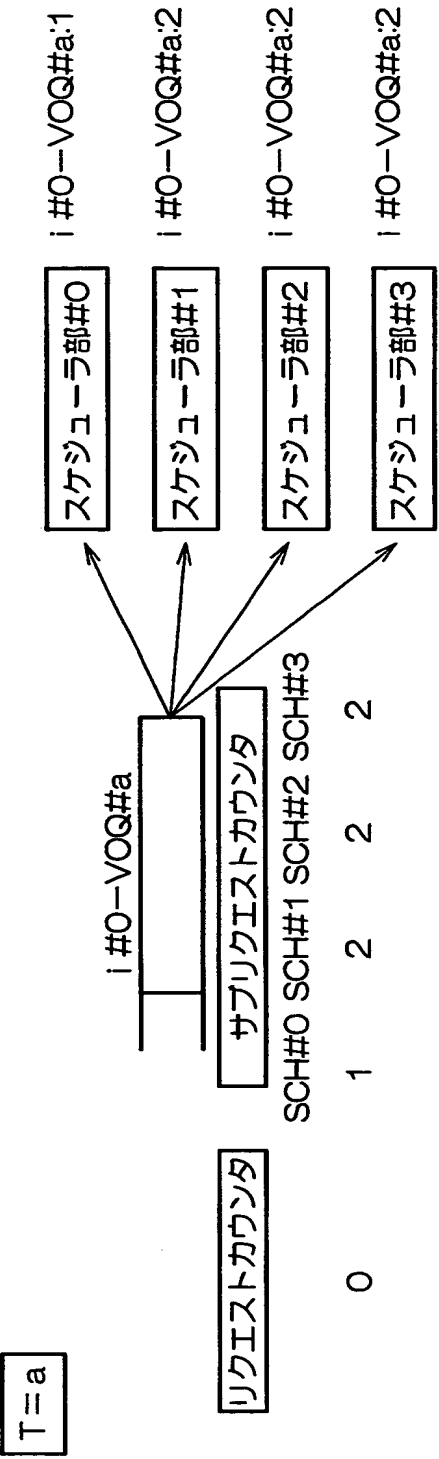


【図 9】



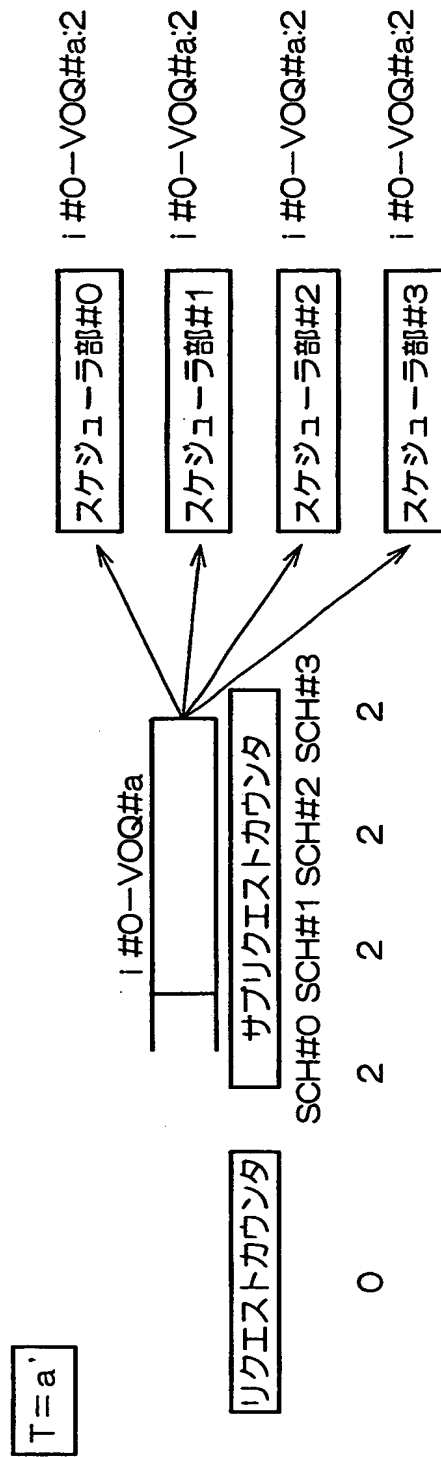
【図 1 0】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図



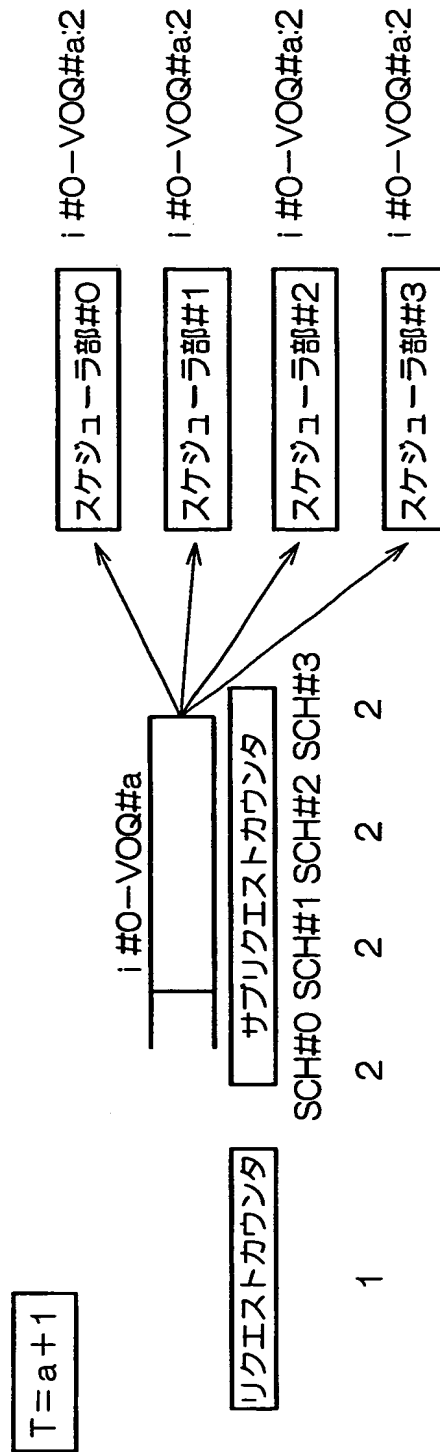
【図 1 1】

リクエスト数の管理機能を入カバッファ部とスケジューラ部に分散配備した場合の具体例を示す図



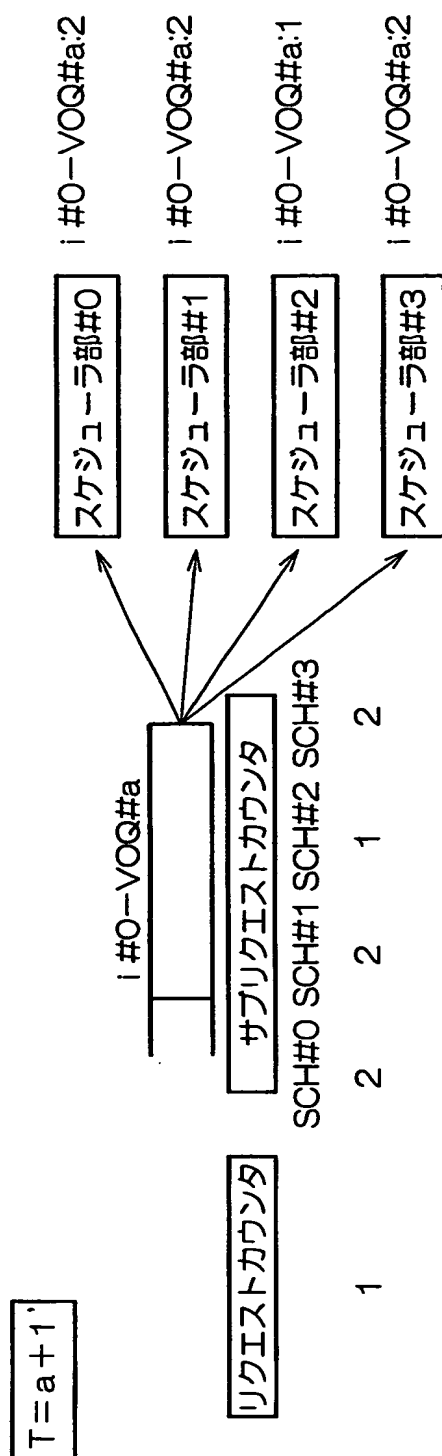
【図 1 2】

リクエスト数の管理機能を入カバッファ部とスケジューラ部に
分散配備した場合の具体例を示す図



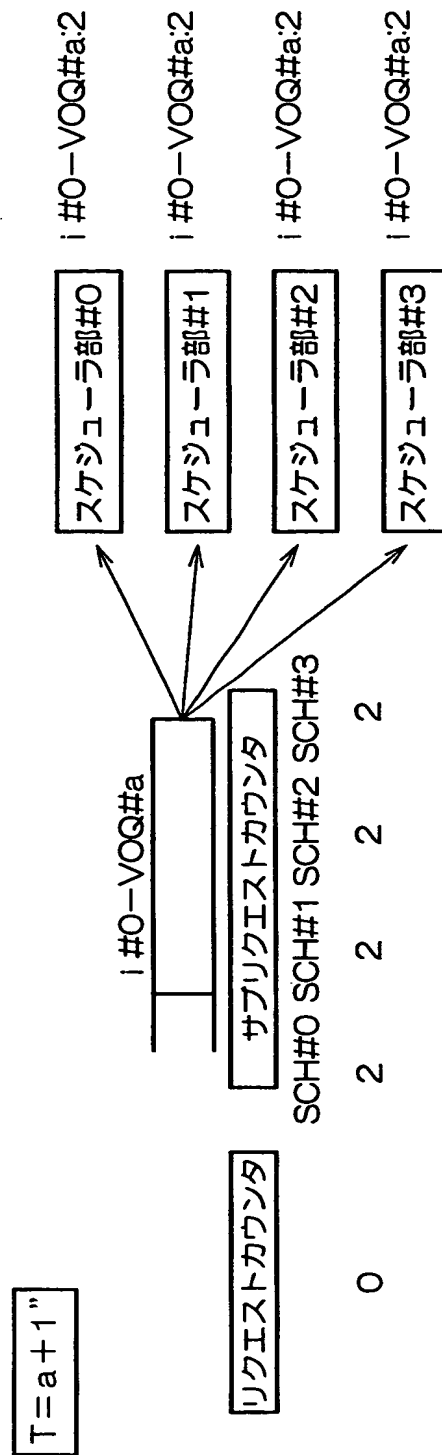
【图 13】

リクエス ト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図



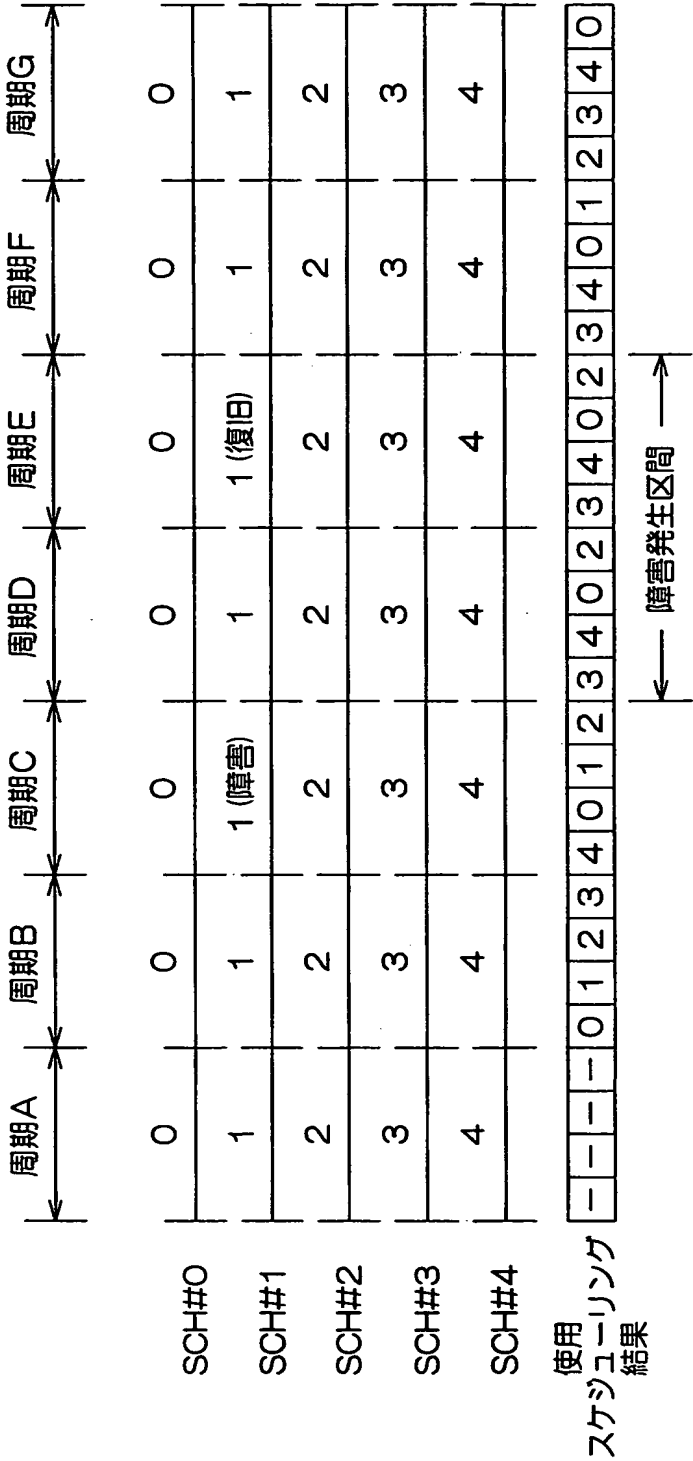
【図 1 4】

リクエスト数の管理機能を入力バッファ部とスケジューラ部に分散配備した場合の具体例を示す図



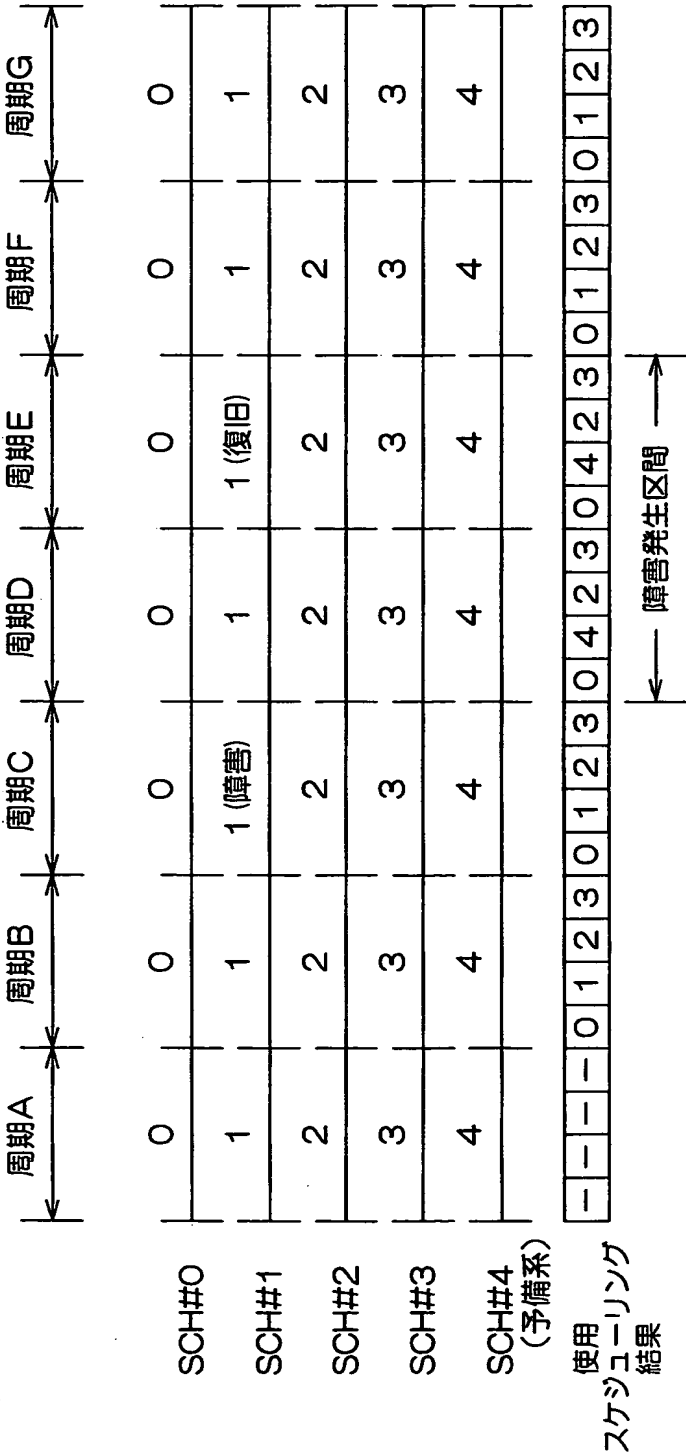
【図 1 5】

冗長系のスケジューラ部を非固定にした場合の具体的な例を示す図



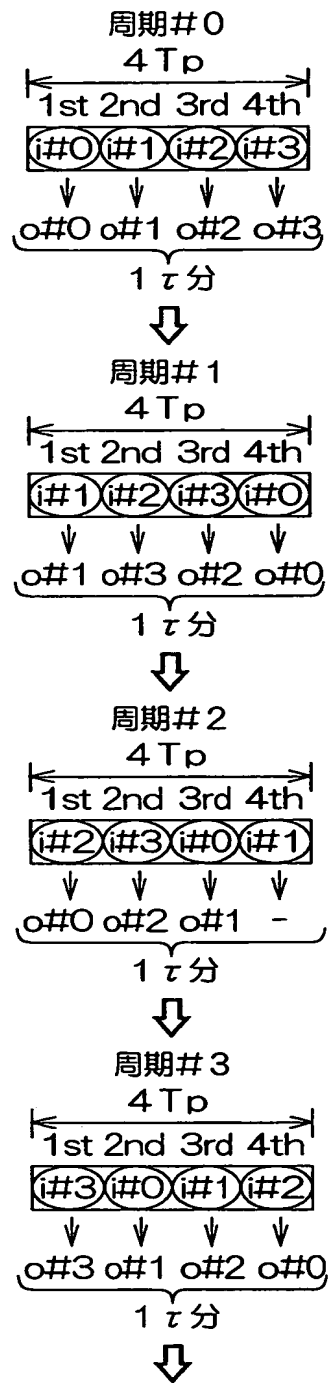
【図 1 6】

冗長系のスケジュール部を固定にした場合の具体的な例を示す図



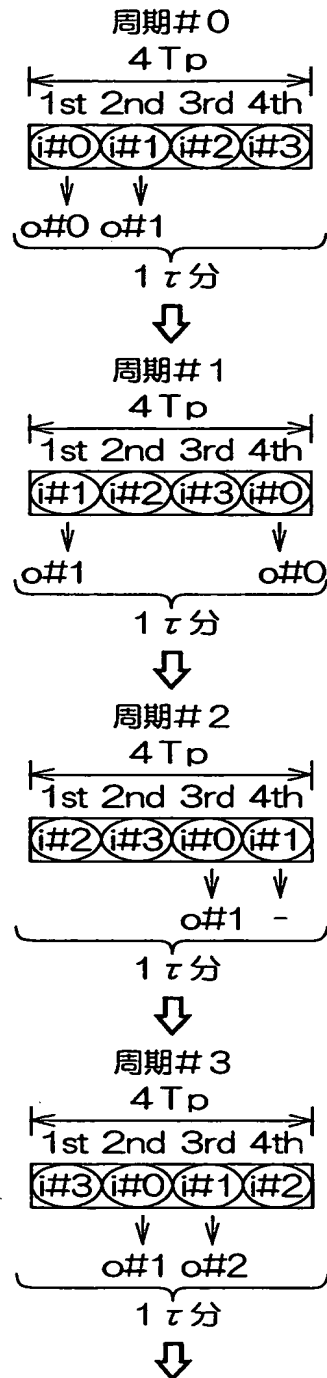
【図 1 7】

4×4のパケットスイッチを構成した場合の
スケジューリング処理の具体例を示す図



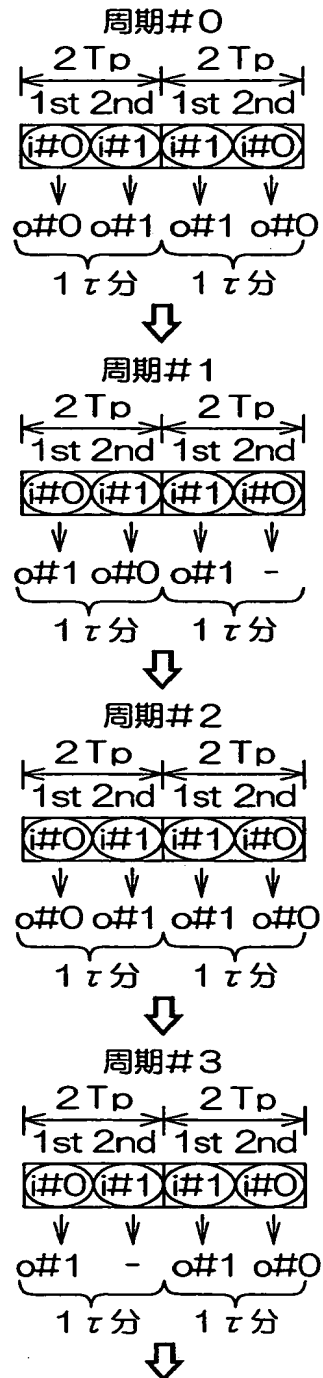
【図 1 8】

2×2の packets スイッチを構成した場合のスケジューリング
処理の具体例を示す図（スループットが低下する場合）



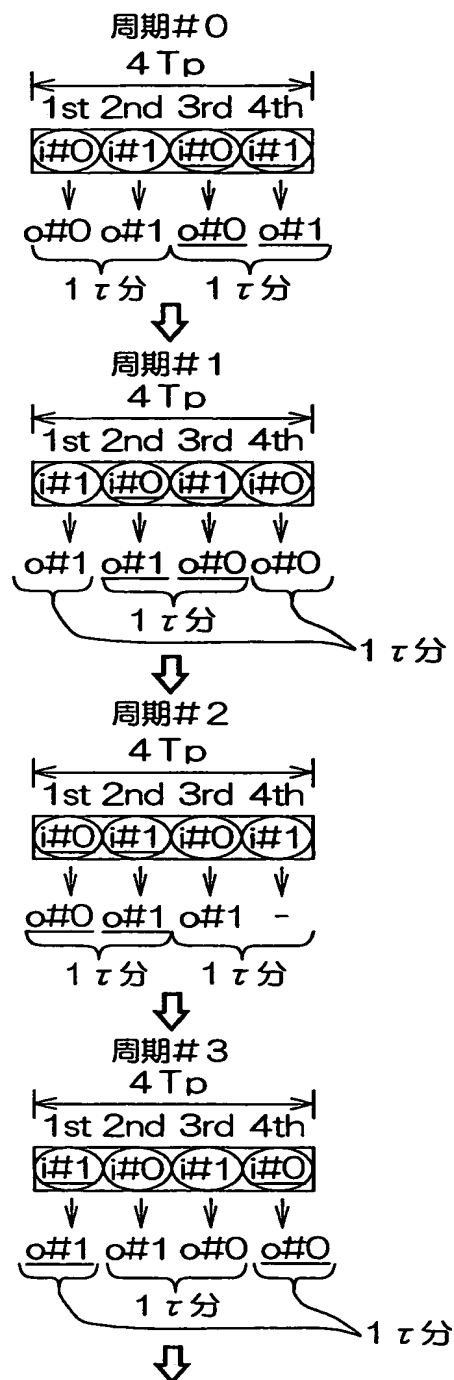
【図 1 9】

2×2の packets スイッチを構成した場合のスケジューリング
処理の具体例を示す図（スループットが低下しない場合）



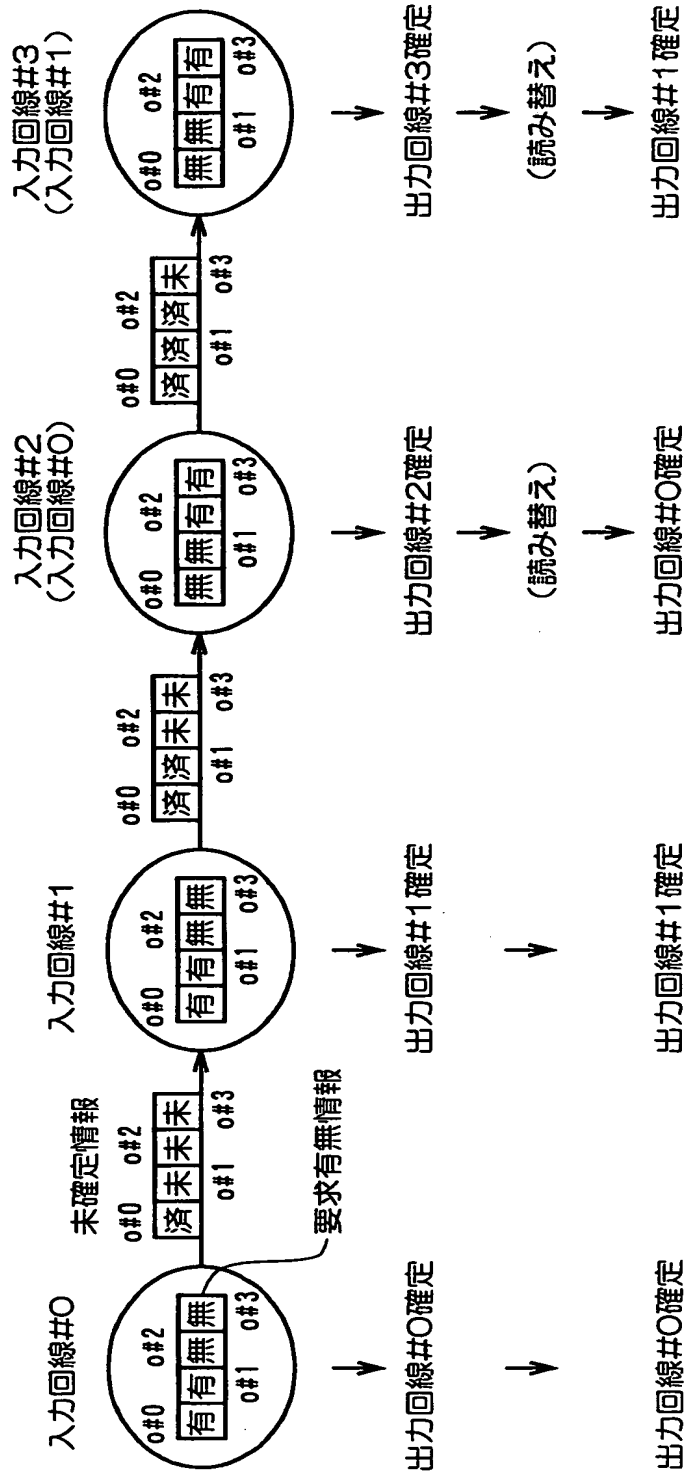
【図 2 0】

2×2のパケットスイッチを構成した場合のスケジューリング
処理の具体例を示す図（スループットが低下しない場合）



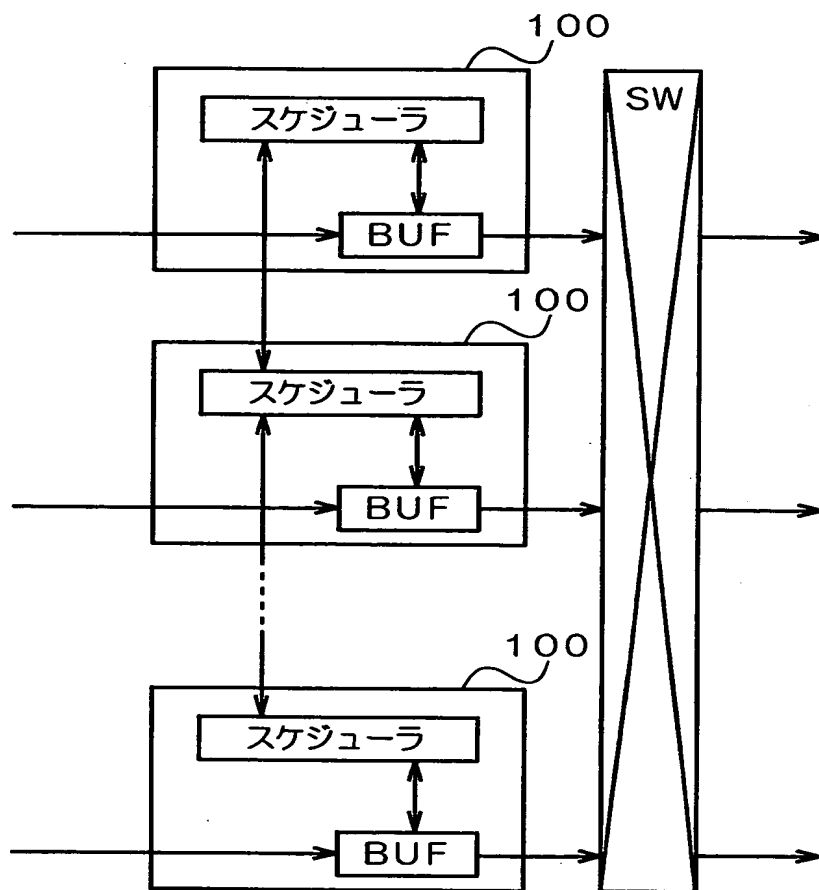
【図 2 1】

1 周期の中で複数のスケジューリング結果を導き出すスケジューリング処理の具体例を示す図



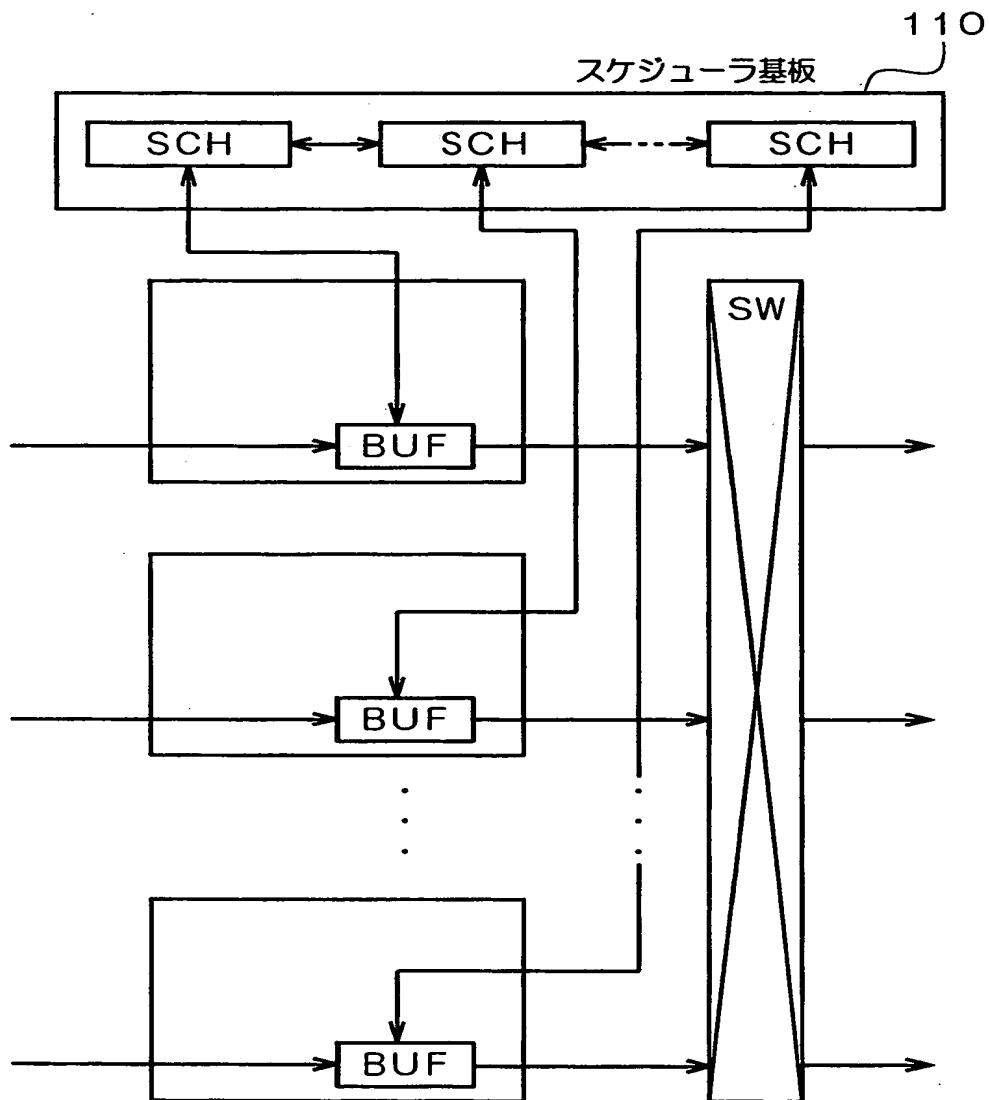
【図 2 2】

スケジューリング機能を分散して配備した
従来のパケットスイッチの構成図



【図 23】

スケジューリング機能を集中して配備した
従来のパケットスイッチの構成図



【書類名】 要約書

【要約】

【課題】 実装時の制約が少なく、無駄な構成を低減でき、しかも拡張性を有するパケットスイッチを提供すること。

【解決手段】 パケットスイッチ 1 0 0 は、スイッチ部 (SW) 1 0、N 個の入力バッファ部 2 0、 α 個のスケジューラ部 3 0 を含んで構成されている。 α 個のスケジューラ部 3 0 によって並列に、しかも互いに独立にスケジューリング処理が行われる。各入力バッファ部 2 0 は、 α 個のスケジューラ部 3 0 のスケジューリング結果を巡回的に利用する。

【選択図】 図 1

出 願 人 履 歴 情 報

識別番号

[000005223]

1. 変更年月日

1996年 3月26日

[変更理由]

住所変更

住 所

神奈川県川崎市中原区上小田中4丁目1番1号

氏 名

富士通株式会社